

# Exploring the Limits of Graph Invariant- and Spectrum-Based Discrimination of (Sub)structures<sup>‡</sup>

Christoph Rücker,\* Gerta Rücker, and Markus Meringer

Department of Mathematics, Universität Bayreuth, D-95440 Bayreuth, Germany

Received December 1, 2001

The limits of a recently proposed computer method for finding all distinct substructures of a chemical structure are systematically explored within comprehensive graph samples which serve as supersets of the graphs corresponding to saturated hydrocarbons, both acyclic (up to  $n = 20$ ) and (poly)cyclic (up to  $n = 10$ ). Several pairs of smallest graphs and compounds are identified that cannot be distinguished using selected combinations of invariants such as combinations of Balaban's index  $J$  and graph matrix eigenvalues. As the most important result, it can now be stated that the computer program NIMSG, using  $J$  and distance eigenvalues, is safe within the domain of mono- through tetracyclic saturated hydrocarbon substructures up to  $n = 10$  (oligocyclic decanes) and of all acyclic alkane substructures up to  $n = 19$  (nonadecanes), i.e., it will not miss any of these substructures. For the regions surrounding this safe domain, upper limits are found for the numbers of substructures that may be lost in the worst case, and these are low. This taken together means that the computer program can be reasonably employed in chemistry whenever one is interested in finding the saturated hydrocarbon substructures. As to unsaturated and heteroatom containing substructures, there are reasons to conjecture that the method's resolving power for them is similar.

## INTRODUCTION

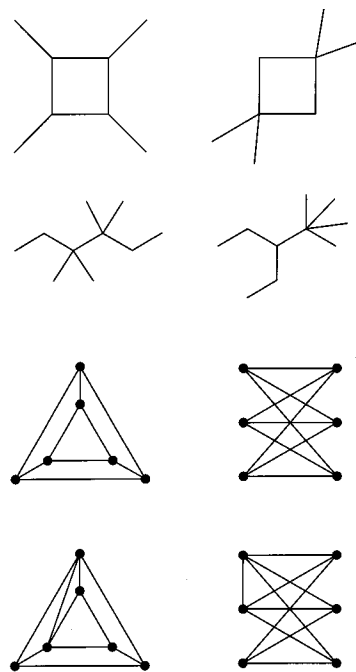
Substructures and subgraphs of chemical structures are becoming increasingly important in description of chemical compounds' properties and reactivity,<sup>1a</sup> in similarity and complexity considerations,<sup>1b,c</sup> in physical and biological property prediction,<sup>1d</sup> and in automatic structure elucidation from spectral data.<sup>2</sup> We recently developed computer programs capable of finding all connected subgraphs in simple graphs,<sup>3</sup> all connected substructures and distinct connected substructures in colored multigraphs and chemical structures,<sup>4</sup> and all connected substructures and subgraphs and distinct connected substructures and subgraphs in colored multigraphs and chemical structures.<sup>5</sup> In such an endeavor the ability to distinguish very similar graphs is obviously a central issue and is in fact the limiting factor. Since a fast computer method for reliably discriminating all nonisomorphic graphs was not at our hands, the best we could do was to use graph invariants of discriminating power as high as possible.

A graph herein is understood to be simple, connected, and undirected. It contains  $n$  vertices,  $m$  edges, and  $c = m - n + 1$  cycles. A graph invariant is a number calculated for a graph from its structure according to a well-defined procedure, its value is independent of how the graph is drawn or how its vertices are numbered. Being a simple number, a graph invariant carries less information than the graph itself, and this loss of information results in graph invariants being more or less degenerate, i.e., nonisomorphic graphs may have the same value of a particular invariant.

An easy-to-calculate graph invariant which is nevertheless considered rather well-discriminating is Balaban's index  $J$ .<sup>6</sup> Index  $J$  is of low degeneracy (has high discriminating power)

\* Corresponding author phone: +49-921-553386; fax: +49-921-553385; e-mail: Christoph.Ruecker@uni-bayreuth.de.

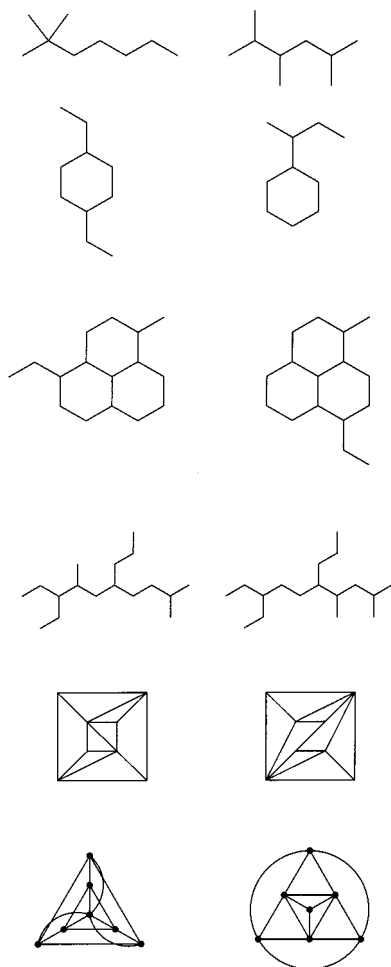
<sup>‡</sup> Dedicated to Professor A. T. Balaban on the occasion of his 71st birthday.



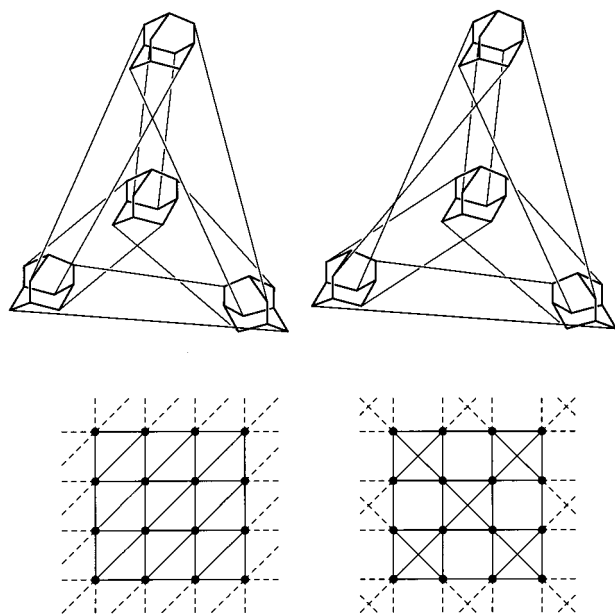
**Figure 1.** Pairs of  $J$ -equivalent but not isospectral graphs.

compared to several other well-known graph invariants, in that the smallest  $J$ -equivalent simple graphs have  $n = 6$  vertices, the smallest  $J$ -equivalent tree (=acyclic) graphs are found in the  $n = 10$  family, and the smallest  $J$ -equivalent alkanes (4-trees) are dodecanes.<sup>7</sup>

A better resolution should be achievable by using, instead of one graph invariant, a sequence of several graph invariants,<sup>8</sup> such as a spectrum. A graph's spectrum is the sequence of the eigenvalues of its adjacency matrix, a one-dimensional array of  $n$  graph invariants. The spectrum still contains less information than the graph itself, i.e., two nonisomorphic graphs may exhibit the same spectrum, in which case they



**Figure 2.** Pairs of isospectral but not  $J$ -equivalent graphs. The last four pairs are even distance-isospectral.



**Figure 3.** Pairs of graphs that are both  $J$ -equivalent and isospectral (and distance-isospectral).

are called isospectral or cospectral graphs. The smallest isospectral connected simple graphs have  $n = 6$  vertices,<sup>9,10</sup> the smallest isospectral tree graphs are in the  $n = 8$  family,<sup>9</sup> and the smallest isospectral alkanes are nonanes.<sup>11</sup> These numerical results, when compared to those for index  $J$  cited

**Table 1.** Numbers of Graphs with  $n$  Vertices and  $m$  Edges for  $n \leq 10$

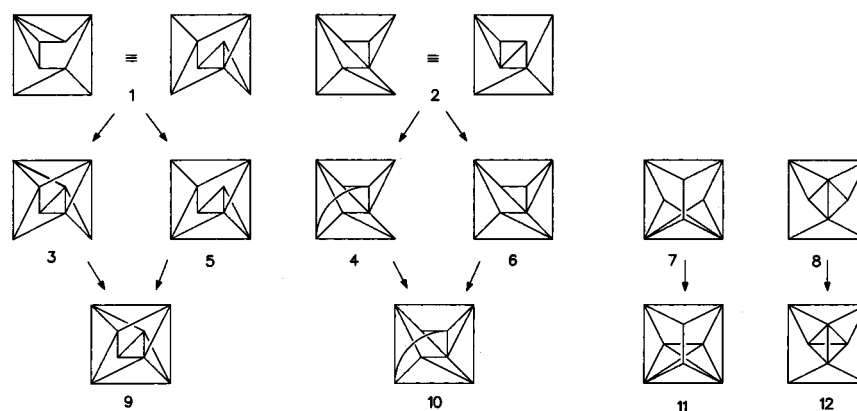
$n$	1	2	3	4	5	6	7	8	9	10
$m$	1									
0	1									
1		1								
2			1							
3				1						
4					2					
5						3				
6							5			
7								6		
8									11	
9										23
10										
11										
12										
13										
14										
15										
16										
17										
18										
19										
20										
21										
22										
23										
24										
25										
26										
27										
28										
29										
30										
31										
32										
33										
34										
35										
36										
37										
38										
39										
40										
41										
42										
43										
44										
45										
$\Sigma$	1	1	2	6	21	112	853	11117	261080	11716571

above, are somewhat unexpected, they emphasize the extraordinarily high resolving power of  $J$ . Whether or not index  $J$  is generally better resolving than the spectrum was never investigated. It was, however, proven that for increasing  $n$  the fraction of isospectral trees among all trees approaches 1, i.e., "almost all trees are cospectral".<sup>12</sup>

More discriminant than the usual (adjacency matrix) spectrum seems to be the graph distance spectrum, i.e., the sequence of eigenvalues of the graph distance matrix.<sup>13</sup> The smallest distance-isospectral trees have  $n = 17$  vertices and are alkane (heptadecane) graphs,<sup>14,15</sup> while the smallest distance-isospectral simple graphs were not known at the beginning of this study. So neither simple-number graph invariants nor spectra seem to uniquely characterize a graph, i.e., discriminate it from all nonisomorphic graphs.

We had found that as a rule of thumb pairs of  $J$ -equivalent graphs are discriminated by their adjacency or distance spectra (see Figure 1), and conversely typical isospectral and even distance-isospectral graphs are discriminated by their  $J$  values (see Figure 2). So we formulated the working hypothesis that this will be generally the case, at least for small and not too complex (molecular) graphs. Accordingly, we decided to use for graph discrimination in our computer program NIMSG the combination of  $J$  and adjacency or distance eigenvalues.<sup>4</sup>

Of course a pair of graphs that are at the same time  $J$ -equivalent and isospectral cannot be distinguished by this method. Thus if two such graphs appear both as subgraphs



**Figure 4.** Smallest graphs that are pairwise both  $J$ -equivalent and isospectral (and moreover distance-isospectral).

**Table 2.** Numbers of Distinct  $J$  Values and Resolution by  $J$  for Graphs of  $n \leq 10$

	n	1	2	3	4	5	6	7	8	9	10				
m															
0	1														
1		1													
2			1												
3				1											
4					2										
5						3									
6							6	1							
7															
8															
9															
10															
11															
12															
13															
14															
15															
16															
17															
18															
19															
20															
21															
22															
23															
24															
25															
26															
27															
28															
29															
30															
31															
32															
33															
34															
35															
36															
37															
38															
39															
40															
41															
42															
43															
44															
45															
	1	1	2	6	21	107	0.955	762	0.893	8200	0.738	138749	0.531	3648987	0.311

in a graph, the result will be a wrong (low by 1) number of distinct subgraphs. Before the present study was initiated we knew of only a few such pairs of graphs, e.g. two regular cubic (degree of each vertex equals 3) simple graphs of 40 vertices<sup>16</sup> or two nonmolecular graphs of 16 vertices (all vertex degrees equal to 6).<sup>17</sup> These graphs are shown in Figure 3, and further examples can be found in the work of Weisfeiler<sup>18</sup> and Mathon.<sup>19</sup> In the context of molecular structures all these graphs seemed irrelevant, most for their high vertex degrees ( $>4$ ), the first-mentioned pair for their

size in combination with their regularity and the unfavorable geometry of any 3D-realization.

Treating the complete graphs  $K_n$  up to  $n = 7$  as tests of our program NIMSG had resulted in the correct numbers of distinct subgraphs.<sup>20</sup> So we knew that at least up to and including  $n = 7$  no such "dangerous" pairs of simultaneously  $J$ -equivalent and isospectral simple graphs exist. The aim of the present work was to find out whether such dangerous pairs are a realistic threat in finding molecular substructures, in particular in the application of NIMSG to molecular struc-

**Table 3.** Numbers of Distinct Adjacency Spectra and *Resolution* by Adjacency Spectra for Graphs of  $n \leq 10$ 

n	1	2	3	4	5	6	7	8	9	10				
m														
0	1													
1		1												
2			1											
3				1										
4					2									
5						3								
6							6	1						
7									11	1				
8														
9														
10														
11														
12														
13														
14														
15														
16														
17														
18														
19														
20														
21														
22														
23														
24														
25														
26														
27														
28														
29														
30														
31														
32														
33														
34														
35														
36														
37														
38														
39														
40														
41														
42														
43														
44														
45														
1	1	2	6	21	111	0.991	821	0.962	10423	0.938	236064	0.904	10375797	0.884

tures. This was to be done by systematically identifying the smallest such pairs of graphs.

## RESULTS AND DISCUSSION

So questions arose as to the size and identities of the smallest simple graphs simultaneously being  $J$ -equivalent and isospectral and to the nature of such graphs — “molecular” or not. Unfortunately, no simple definition of a molecular graph is available. Therefore in the following we treat the sets of connected simple graphs, of connected simple 4-graphs, of trees and of 4-trees up to a certain vertex number, each of which is a superset of cyclic or acyclic saturated hydrocarbon graphs, respectively.

New hardware now allowed us to fully treat the complete graph  $K_8$ . As it happened, the number of distinct connected subgraphs of  $n = 8$  found was low, 11111 instead of 11117,<sup>20</sup> even if all eight adjacency matrix eigenvalues or all eight distance matrix eigenvalues were used together with  $J$  for discrimination (instead of the routinely employed two adjacency or two distance eigenvalues). In detail, our procedure found 1578 instead of 1579 distinct connected simple graphs of  $n = 8$ ,  $m = 14$  (corresponding to heptacyclic octanes), 1512 instead of 1515 distinct connected simple graphs of  $n = 8$ ,  $m = 15$  (octacyclic octanes), and 1288 instead of 1290 distinct connected simple graphs of  $n$

$= 8$ ,  $m = 16$  (nonacyclic octanes). Each distinct subgraph found occurs in many copies within  $K_8$  due to its high symmetry, e.g. a typical occurrence number of  $n = 8$ ,  $m = 16$  subgraphs in  $K_8$  is 23040. For  $n = 8$ , all other  $m$  ( $7 \leq m \leq 13$  and  $17 \leq m \leq 28$ ), the numbers of distinct connected simple graphs found were correct.<sup>20</sup> At this stage we knew that there must exist a few pairs of graphs with the sought-after combination of properties for  $n = 8$ ,  $m = 14$ –16, but so far it was impossible to identify them. Comfortably, it was also clear that hepta-, octa-, and nonacyclic graphs of eight vertices are not molecular graphs.<sup>21</sup>

The key to the successful identification reported here is a complete generation free of redundancy of all connected simple graphs of  $n = 8$ ,  $m = 14$ , 15, 16, that was now performed using MOLGEN 4.0.<sup>22</sup> Within MOLGEN 4.0, isomorphic graphs are identified, and nonisomorphic graphs are distinguished by a canonical numbering scheme. Calculation of  $J$  and the eigenvalues for all 1579 graphs of  $n = 8$ ,  $m = 14$ , 1515 graphs of  $n = 8$ ,  $m = 15$ , and 1290 graphs of  $n = 8$ ,  $m = 16$  and sorting by  $J$  or/and the eigenvalues within each class led to the following observations:<sup>23</sup>

(i) There are many pairs, triplets, and higher tuples of  $J$ -equivalent graphs in each of these classes.<sup>24</sup>

(ii) There are many pairs and several triplets of isospectral and even distance isospectral graphs in these graph classes.



**Table 4.** Numbers of Distinct Distance Spectra and *Resolution* by Distance Spectra for Graphs of  $n \leq 10$ 

n	1	2	3	4	5	6	7	8	9	10
m										
0	1									
1		1								
2			1							
3				1						
4					2					
5						3				
6							5			
7								11	1	
8										23
9										
10										
11										
12										
13										
14										
15										
16										
17										
18										
19										
20										
21										
22										
23										
24										
25										
26										
27										
28										
29										
30										
31										
32										
33										
34										
35										
36										
37										
38										
39										
40										
41										
42										
43										
44										
45										
1	1	1	2	6	21	112	842	0.987	10785	0.970
2										
3										
4										
5										
6										
7										
8										
9										
10										
11										
12										
13										
14										
15										
16										
17										
18										
19										
20										
21										
22										
23										
24										
25										
26										
27										
28										
29										
30										
31										
32										
33										
34										
35										
36										
37										
38										
39										
40										
41										
42										
43										
44										
45										
10975530										
0.937										

(iii) For  $n = 8$ , there are exactly the following six pairs of graphs which are simultaneously  $J$ -equivalent and isospectral:

$m = 14$  class (heptacyclic octanes): **1** and **2** shown in Figure 4;

$m = 15$  class (octacyclic octanes): **3** and **4**, **5** and **6**, and **7** and **8**;

$m = 16$  class (nonacyclic octanes): **9** and **10** and **11** and **12**.

Four of these are planar graphs (**1,2,6,8**), the others are nonplanar.

Surprisingly, within each such pair of graphs even the distance matrix eigenvalues coincide, i.e., these graphs are pairwise not only  $J$ -equivalent and isospectral, but even distance-isospectral.<sup>25</sup> Furthermore, the Wiener index  $W$  and Hosoya index  $Z$  values (and their building blocks  $p(G,k)$ ) pairwise coincide. In fact these graphs are “topological twin graphs” in the sense of Hosoya,<sup>26</sup> but being  $J$ -equivalent they are even more similar to one another than required by the definition of topological twins.<sup>27</sup> Furthermore, with respect to the number of edges graphs **1** and **2** are smaller than Hosoya’s smallest topological twins.

These graphs are genetically related, they form two families: From **1** both **3** and **5** can be formed by addition of an edge and adding the respective other edge to either of these results in **9**. Likewise, from **2** (the twin of **1**) by adding an edge **4** and **6** (the twins of **3** and **5**) can be formed, and

either of these leads to **10** (the twin of **9**) by adding the respective other edge. The second family is formed by **7**, its twin **8**, and **11** and its twin **12**, where the latter result from addition of an edge to either of the former. As anticipated, all these graphs are nonmolecular graphs due to vertex degrees exceeding 4.

Thus a partial answer to the question on the limits of validity of our working hypothesis above was found. However, from these findings the following questions arose: Are distance spectra really more discriminating than adjacency spectra? How frequent are  $J$ -equivalent graphs, isospectral graphs, distance-isospectral graphs, graphs both  $J$ -equivalent and isospectral, graphs both  $J$ -equivalent and distance-isospectral among the graphs of  $n > 8$ , and in particular among the molecular graphs of that size? What are the smallest molecular graphs simultaneously  $J$ -equivalent and isospectral/distance-isospectral?

**Graphs of  $n \leq 10$ .** To get an idea on possible answers to these questions we decided first to systematically look for degeneracies in  $J$  and adjacency and distance spectra within the set of all connected simple graphs of up to  $n = 10$ , which using MOLGEN 4.0 seemed to be a realistic task.

Thus all connected simple graphs of 1, 2, ..., 10 vertices (nearly 12 million graphs) were generated using MOLGEN 4.0 in classes of constant numbers of vertices and edges, their  $J$  values and adjacency and distance spectra were calculated, the numbers of distinct  $J$  values and distinct

**Table 5.** Numbers of Graphs with Distinct Combination of  $J$  and Adjacency Spectrum and *Resolution* by This Combination for  $n \leq 10$ 

	n	1	2	3	4	5	6	7	8	9	10
m	0	1									
1			1								
2				1							
3					2						
4					2						
5					1						
6					1						
7						3					
8						5	6				
9						5	13	11			
10						4	19	33			
11						2	22	67	23	1	
12						1	20	107	236	1	47
13							14	132	486	1	240
14							9	138	814	1	797
15							5	126	1169	1	2075
16							2	95	1454	1	4494
17							1	64	1578	0.999	13849
18							1	40	1512	0.998	20282
19								21	1288	0.998	26566
20								10	970	1	31268
21								5	658	1	33163
22								2	400	1	31727
23								1	220	1	27505
24								1	114	1	21647
25									56	1	15442
26									24	1	10036
27									11	1	5957
28									5	1	3238
29									2	1	1633
30									1	1	770
31									1	1	344
32											148
33											63
34											25
35											11
36											5
37											2
38											1
39											1
40											1
41											1
42											1
43											1
44											1
45											1
		1	1	2	6	21	112	853	11111	0.999	259737
										0.995	11612987
											0.991

adjacency and distance spectra were determined via sorting by  $J$  or the eigenvalues, respectively (two spectra are distinct if they differ in at least one eigenvalue). The results are shown in Tables 1–6. In the tables every fifth row is underlined for better orientation. Table 1 gives the numbers of connected simple graphs in classes of constant  $n$  and  $m$ , as known<sup>20</sup> and as generated by MOLGEN 4.0. These numbers serve as reference values against which to compare the entries in Tables 2–6. Table 2 gives the numbers of distinct  $J$  values within each  $n, m$ -class,<sup>28</sup> Tables 3 and 4 show the numbers of distinct adjacency spectra and distinct distance spectra. Tables 5 and 6 give the numbers of distinct combinations of  $J$  and adjacency spectra and of  $J$  and distance spectra, respectively, for the same classes of graphs. Tables 2–6 also show *in italics* the resolution of the respective graph invariant (combination), i.e., each italic entry is the entry left to it divided by the corresponding entry in Table 1. In the tables the “dangerous” region, the range where the particular invariant (combination) cannot uniquely characterize all graphs, is shaded.

Tables 2–6 all give qualitatively the same picture: The resolutions (discriminating powers) of the graph invariants gradually drop for increasing  $n$ . For increasing  $m$  within each  $n$  the discriminating powers initially drop and then pass through a minimum (printed in bold), finally approaching 1 again. The latter feature is explained by the fact that for increasing  $m$  the numbers of distinct graphs first increase

but then decrease again until the second-highest and highest  $m$  classes contain only one graph each, the  $K_n$ -minus-an-edge and  $K_n$  graphs, so degeneracy in these classes cannot exist.

Huge differences are seen in the discriminating power of the graph invariants considered here:

(i) Index  $J$  is very good for acyclic and oligocyclic graphs (the few first entries in each column), i.e., the domain of real molecular species. In that region  $J$  is even better than the adjacency spectrum. However, down the columns, i.e., for polycyclic graphs,  $J$ 's resolution sharply drops, so that most graphs are better resolved by their spectra.

The different behavior of  $J$  for acyclic and polycyclic graphs may be understood:  $J$  exploits the differences in the (topological) distances between vertices in a graph, more exactly the differences between the distance sums. In going from a tree to a polycycle, long distances are replaced by shorter ones, those that are present in any graph. In the extreme case, the  $K_n$ , all distances are 1 and all distance sums equal  $n - 1$ . So in that direction the distances (and their sums) tend to equalize for the vertices in a graph and between isomeric graphs as well.

(ii) The distance spectrum is always at least as discriminating as the adjacency spectrum.

(iii) The combinations of  $J$  and spectra, particularly  $J$  and the distance spectrum, are unrivalled, as expected.<sup>29</sup> In Table 6 in each column the first five resolution entries are 1, that is, the domain of acyclic to tetracyclic graphs (saturated

**Table 6.** Numbers of Graphs with Distinct Combination of  $J$  and Distance Spectrum and *Resolution* by This Combination for  $n \leq 10$ 

n	1	2	3	4	5	6	7	8	9	10
m										
0	1									
1		1								
2			1							
3				1						
4					2					
5						3				
6							5			
7								11		
8									23	1
9										47
10										
11										
12										
13										
14										
15										
16										
17										
18										
19										
20										
21										
22										
23										
24										
25										
26										
27										
28										
29										
30										
31										
32										
33										
34										
35										
36										
37										
38										
39										
40										
41										
42										
43										
44										
45										
Σ	1	1	2	6	21	112	853	11111	0.999	259910

hydrocarbons) of up to  $n = 10$  is “safe” if the combination of  $J$  and distance spectrum is used for discrimination.

**4-Graphs of  $n \leq 10$ .** Program NIMSG for finding distinct substructures was developed primarily for chemistry, where one is mostly interested in acyclic through oligocyclic graphs (e.g. for  $n \leq 10$ ,  $n - 1 \leq m \leq \sim n + 5$ ) and in particular in graphs of vertex degrees not exceeding 4, the valency of carbon (so-called 4-graphs). The above procedure was therefore repeated for simple connected 4-graphs, the graph sample most closely approximating the acyclic and oligo-through polycyclic saturated hydrocarbons, up to  $n = 10$ . The results are shown in Tables 7–12. Contrary to naïve expectation, the resolution in this sample is not decisively better than in the sample of all graphs, so that the resolution problem essentially remains the same. Though the numbers of 4-graphs (Table 7) are lower than those of all graphs (Table 1), and often far lower, particularly in higher  $m$  classes, the 4-graphs are a more uniform group, so finding differences among them is more difficult. Index  $J$  suffers most from this fact.

The observations made in the sample of all graphs are reproduced in the 4-graphs. From Table 12 it again (and necessarily) follows that all acyclic to tetracyclic saturated hydrocarbons of up to at least  $n = 10$  are distinguished by the combination of  $J$  and distance spectrum.

The pair of smallest  $J$ -equivalent and isospectral 4-graphs was identified in the  $n = 9$ ,  $m = 12$  class (tetracyclic nonanes, Table 11), graphs **13** and **14**; the pair of smallest

**Table 7.** Numbers of 4-Graphs with  $n$  Vertices and  $m$  Edges for  $n \leq 10$ 

n	1	2	3	4	5	6	7	8	9	10
m										
0	1									
1		1								
2			1							
3				1						
4					2					
5						3				
6							5			
7								12		
8									9	
9										29
10										
11										
12										
13										
14										
15										
16										
17										
18										
19										
20										
Σ	1	1	2	6	21	78	353	1929	12207	89402

$J$ -equivalent and distance-isospectral 4-graphs was found in the  $n = 9$ ,  $m = 13$  class (pentacyclic nonanes, Table 12), graphs **15** and **16**. These graphs are shown in Figure 5. Though as 4-graphs they fulfill the formal condition for molecular graphs and though they are planar graphs, a chemist will doubt the viability of their molecular counterparts, due to their presumably extremely strained nature: No reasonable geometric structures (having usual bond lengths,

**Table 8.** Numbers of Distinct  $J$  Values and Resolution by  $J$  for 4-Graphs of  $n \leq 10$ 

n	1	2	3	4	5	6	7	8	9	10					
m	0	1													
0	1														
1		1													
2			1												
3				1											
4					2										
5						3									
6							5								
7								12							
8									1						
9										9					
10															
11															
12															
13															
14															
15															
16															
17															
18															
19															
20															
1	1	1	2	6	21	73	0.936	317	0.898	1662	0.862	10512	0.861	77232	0.864

**Table 9.** Numbers of 4-Graphs with Distinct Adjacency Spectrum and Resolution by Adjacency Spectrum for  $n \leq 10$ 

n	1	2	3	4	5	6	7	8	9	10				
m	0	1												
0	1													
1		1												
2			1											
3				2										
4					3									
5						5								
6							12							
7								9	1					
8										18				
9														
10														
11														
12														
13														
14														
15														
16														
17														
18														
19														
20														
1	1	1	2	6	21	78	345	0.977	1856	0.962	11414	0.935	82824	0.926

**Table 10.** Numbers of 4-Graphs with Distinct Distance Spectrum and Resolution by Distance Spectrum for  $n \leq 10$ 

n	1	2	3	4	5	6	7	8	9	10			
m	0	1											
0	1												
1		1											
2			1										
3				2									
4					3								
5						5							
6							9						
7								18	1				
8													
9													
10													
11													
12													
13													
14													
15													
16													
17													
18													
19													
20													
1	1	1	2	6	21	78	353	1927	0.999	12179	0.998	89240	0.998

bond angles and dihedral angles) will be available to such hypothetical hydrocarbon molecules. Graphs **13** and **14** correspond to substituted tetracyclooctanes of very unusual geometry, more specifically, **13** depicts a bridged [3.2.1]-propellane, **14** a doubly annelated bicyclo[1.1.1]pentane.<sup>30</sup> **15** and **16** correspond to pentacyclononanes, the former to a doubly bridged [3.3.1]propellane, the latter to a bridged

**Table 11.** Numbers of 4-Graphs with Distinct Combination of  $J$  and Adjacency Spectrum and Resolution by This Combination for  $n \leq 10$ 

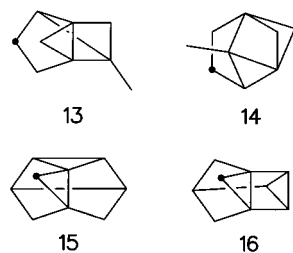
n	1	2	3	4	5	6	7	8	9	10		
m	0	1										
1			1									
2				1								
3					2							
4						3						
5							5					
6								1				
7									9			
8										18		
9												
10												
11												
12												
13												
14												
15												
16												
17												
18												
19												
20												
1	1	1	2	6	21	78	353	1929	12167	0.997	89247	0.998

**Table 12.** Numbers of 4-Graphs with Distinct Combination of  $J$  and Distance Spectrum and Resolution by This Combination for  $n \leq 10$ 

	n	1	2	3	4	5	6	7	8	9	10												
m	0	1																					
1			1																				
2				1																			
3					2																		
4						3																	
5							5																
6								9															
7									18														
8										35	1												
9											75	1											
10													1										
11														1									
12															1								
13																1							
14																	1						
15																		1					
16																			1				
17																				1			
18																					1		
19																						1	
20																							1
	1	1	2	6	21	78	353	1929	12189	0.999	89329	0.999											

[3.2.1]propellane. Not surprisingly, not a single compound containing any of the four polycyclic frameworks of Figure 5 is listed in the Beilstein or the CAS Registry file.





**Figure 5.** Pair of smallest 4-graphs that are both  $J$ -equivalent and isospectral (top) and pair of smallest 4-graphs that are both  $J$ -equivalent and distance-isospectral (and isospectral, bottom). For the meaning of black dots see text.

In the  $n = 10$  domain, the smallest pair of  $J$ -equivalent and isospectral 4-graphs was identified (Table 11) in the  $m = 13$  class (tetracyclic decanes), and the smallest pair of  $J$ -equivalent and distance-isospectral 4-graphs was found to have  $m = 14$  edges (pentacyclic decanes, Table 12). Their structures differ from those shown in Figure 5 only in that they bear an additional vertex attached to the one marked with a dot.

From these smallest examples of simultaneously  $J$ -equivalent and (distance-) isospectral 4-graphs it is concluded that such graphs probably are too complex, too polycyclic for their molecular counterparts to be capable of existence. Other 4-graph pairs of  $n = 9$  or  $10$  being  $J$ -equivalent and (distance-) isospectral have even higher  $m$  values, meaning that molecular counterparts would contain even more cycles than those found above, and therefore will tend to be even more strained. This means that it is reasonably safe to use the combination of  $J$  and (distance) spectrum for identifying distinct molecular substructures and molecular subgraphs, at least in the size range investigated here.

Finally, since NIMSG uses along with  $J$  only two (distance) eigenvalues rather than the complete (distance) spectrum for discrimination among subgraphs, for the sample of 4-graphs we repeated the described procedure, but using only two adjacency eigenvalues or two distance eigenvalues. By systematic variation it was found that the combinations  $\lambda_2$  and  $\lambda_3$  and  $\delta_1$  and  $\delta_n$  (used in two published variants of NIMSG<sup>4</sup>) are not optimal. The most discriminating combinations we were able to find are  $\lambda_3$  and  $\lambda_n$  (the third and the last adjacency eigenvalues) and  $\delta_2$  and  $\delta_{n-1}$  (the second and second-last distance eigenvalues). As a consequence, NIMSG was now improved accordingly. The results shown in Tables 13 and 14 allow for the estimation of the “safety” of the new NIMSG versions or the risk of obtaining too few distinct substructures/subgraphs. As was to be expected, the results in Table 13 are somewhat inferior to those in Table 11, those in Table 14 are inferior to those in Table 12. However, resolution losses due to using only two instead of all eigenvalues appear in the high  $m$  region only, that is for graphs certainly not molecular.

**Tree Graphs of  $n \leq 20$ .** For trees (uppermost entry in each column in Tables 1–14) the resolution of the combinations  $J$  and adjacency spectrum and  $J$  and distance spectrum is perfect in our graph sample of  $n \leq 10$ , as expected (recall that the first degeneracy of the adjacency spectrum and of the distance spectrum for trees are known to occur for  $n = 8$  and  $n = 17$ , respectively). To fathom corresponding limits we additionally generated all trees of up to  $n = 20$  and searched their  $J$  values and spectra.<sup>31</sup> The results are given

**Table 13.** Numbers of 4-Graphs with Distinct Combination of  $J$ ,  $\lambda_3$ , and  $\lambda_n$  and Resolution by This Combination for  $n \leq 10$

$n$	1	2	3	4	5	6	7	8	9	10
$m$	0	1								
1	1									
2		1								
3			1							
4				2						
5					3					
6						5				
7							12			
8								9		
9									18	
10										1
11										
12										
13										
14										
15										
16										
17										
18										
19										
20										
1	1	1	2	6	21	78	353	1928	0.999	12158
2									0.996	89235
3										0.998

**Table 14.** Numbers of 4-Graphs with Distinct Combination of  $J$ ,  $\delta_2$ , and  $\delta_{n-1}$  and Resolution by This Combination for  $n \leq 10$

$n$	1	2	3	4	5	6	7	8	9	10
$m$	0	1								
1	1									
2		1								
3			1							
4				2						
5					3					
6						5				
7							12			
8								9		
9									18	
10										1
11										
12										
13										
14										
15										
16										
17										
18										
19										
20										
1	1	1	2	6	21	78	353	1928	0.999	12179
2									0.998	89309
3										0.999

in Table 15. Where differences are found between the resolutions of  $J$  and the spectra for tree graphs, single index  $J$  is more discriminating than the complete adjacency spectrum but less discriminating than the distance spectrum. First degeneracies of both  $J$ /spectrum combinations are encountered for  $n = 20$ , there are two pairs of  $J$ -equivalent and isospectral such trees (17/18 and 19/20 in Figure 6), and of these one pair (19/20) is even distance isospectral. All these trees are 4-trees, i.e., alkanes, eicosanes. In both pairs the structures differ in a position exchange of ethyl and gem-dimethyl substituents, as was discussed earlier.<sup>7a</sup>

**4-Trees of  $n \leq 20$ .** Results for all alkanes  $C_nH_{2n+2}$  of up to  $n = 20$ , generated using MOLGEN 4.0, are given in Table 16. Here as for the general trees  $J$  is more discriminant than the adjacency spectrum but less than the distance spectrum. Within the alkanes the resolution of  $J$  is somewhat higher, that of the adjacency spectrum is somewhat lower than within all trees. It was also checked (not shown in the table) that use of  $\lambda_3$  and  $\lambda_n$  instead of all adjacency eigenvalues and of  $\delta_2$  and  $\delta_{n-1}$  instead of all distance eigenvalues (the NIMSG procedures) does not compromise the complete discrimination among alkanes of up to  $n = 19$  (the nonadecanes).

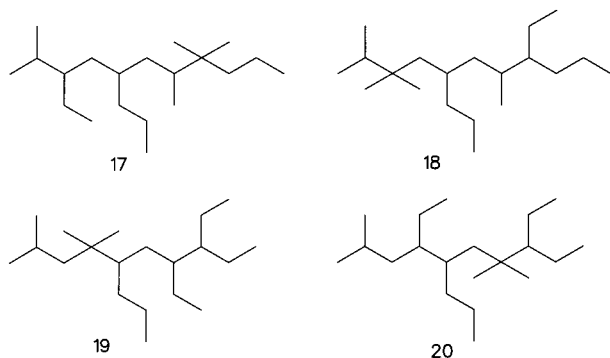
**Concluding Remarks.** Let us emphasize here once more that structure discrimination by combinations of graph invariants (as done in NIMSG) seems to be a simple but

**Table 15.** Numbers of Distinct Values, and Resolution of  $J$ , Adjacency Spectrum, Distance Spectrum, Combination of  $J$  and Adjacency Spectrum, and Combination of  $J$  and Distance Spectrum, for Trees of  $n \leq 20$ 

$n$	#	$J$ #	$J$ res	adj spectrum #	adj spectrum res	dist spectrum #	dist spectrum res	$J$ /adj sp #	$J$ /adj sp res	$J$ /dist sp #	$J$ /dist sp res
1	1	1	1	1	1	1	1	1	1	1	1
2	1	1	1	1	1	1	1	1	1	1	1
3	1	1	1	1	1	1	1	1	1	1	1
4	2	2	1	2	1	2	1	2	1	2	1
5	3	3	1	3	1	3	1	3	1	3	1
6	6	6	1	6	1	6	1	6	1	6	1
7	11	11	1	11	1	11	1	11	1	11	1
8	23	23	1	22	0.957	23	1	23	1	23	1
9	47	47	1	42	0.894	47	1	47	1	47	1
10	106	105	0.991	102	0.962	106	1	106	1	106	1
11	235	234	0.996	204	0.868	235	1	235	1	235	1
12	551	537	0.975	488	0.886	551	1	551	1	551	1
13	1301	1290	0.992	1078	0.829	1301	1	1301	1	1301	1
14	3159	3026	0.958	2723	0.862	3159	1	3159	1	3159	1
15	7741	7609	0.983	6403	0.827	7741	1	7741	1	7741	1
16	19320	18158	0.940	16479	0.853	19320	1	19320	1	19320	1
17	48629	47480	0.976	40313	0.829	48628	0.999	48629	1	48629	1
18	123867	114600	0.925	106135	0.857	123865	0.999	123867	1	123867	1
19	317955	308063	0.969	271295	0.853	317949	0.999	317955	1	317955	1
20	823065	749284	0.910	724455	0.880	823051	0.999	823063	0.999	823064	0.999

**Table 16.** Numbers of Distinct Values, and Resolution of  $J$ , Adjacency Spectrum, Distance Spectrum, Combination of  $J$  and Adjacency Spectrum, and Combination of  $J$  and Distance Spectrum, for 4-Trees (Alkanes  $C_nH_{2n+2}$ ) of  $n \leq 20$ 

$n$	#	$J$ #	$J$ res	adj spectrum #	adj spectrum res	dist spectrum #	dist spectrum res	$J$ /adj sp #	$J$ /adj sp res	$J$ /dist sp #	$J$ /dist sp res
1	1	1	1	1	1	1	1	1	1	1	1
2	1	1	1	1	1	1	1	1	1	1	1
3	1	1	1	1	1	1	1	1	1	1	1
4	2	2	1	2	1	2	1	2	1	2	1
5	3	3	1	3	1	3	1	3	1	3	1
6	5	5	1	5	1	5	1	5	1	5	1
7	9	9	1	9	1	9	1	9	1	9	1
8	18	18	1	18	1	18	1	18	1	18	1
9	35	35	1	30	0.857	35	1	35	1	35	1
10	75	75	1	73	0.973	75	1	75	1	75	1
11	159	159	1	136	0.855	159	1	159	1	159	1
12	355	349	0.983	307	0.865	355	1	355	1	355	1
13	802	799	0.996	652	0.813	802	1	802	1	802	1
14	1858	1808	0.973	1580	0.850	1858	1	1858	1	1858	1
15	4347	4305	0.990	3484	0.801	4347	1	4347	1	4347	1
16	10359	9923	0.958	8573	0.828	10359	1	10359	1	10359	1
17	24894	24516	0.985	19786	0.795	24893	0.999	24894	1	24894	1
18	60523	57331	0.947	50340	0.832	60521	0.999	60523	1	60523	1
19	148284	145206	0.979	122453	0.826	148279	0.999	148284	1	148284	1
20	366319	342886	0.936	313498	0.856	366308	0.999	366317	0.999	366318	0.999

**Figure 6.** Two pairs of both  $J$ -equivalent and isospectral eicosanes. The bottom pair is distance-isospectral as well. These are the smallest alkane graphs having these properties.

only approximate solution to a difficult problem. Here we considered graphs corresponding to a superset of saturated hydrocarbons (of rather low carbon number) only, so we cannot say anything about the discrimination of real chemical structures other than saturated hydrocarbons. Most molecular structures, containing multiple bonds and heteroatoms, are to be represented by colored multigraphs. Obviously, there are many more colored multigraphs than simple graphs for each vertex number  $n$ , so that their discrimination seems to be even more difficult. On the other hand, we carefully

included information on multiple bonds and heteroatoms into  $J$  and the spectra used in NIMSG,<sup>4,32</sup> hopefully raising the discriminating power of the procedure to a level sufficient for practical purposes. Further, in mathematical graph theory experience is that the graph isomorphism problem is more difficult for simple graphs than for colored multigraphs, the former lacking distinguishing features. To test this point would require one to have a comprehensive sample of molecular colored multigraphs, which obviously is not at hand for any  $n$ .

After proving that “almost all trees are cospectral”, Schwenk raised the question whether the same is true for almost all graphs.<sup>12</sup> He did, however, not answer this question, nor did he give a conjecture, due to considerable differences in the mathematical properties of trees on one side and (general) graphs on the other. We here obtained at least some experimental information relevant to this issue. In the adjacency spectrum column in Table 15, resolution values oscillate and only slowly decrease for increasing  $n$ , so that one would probably not have predicted Schwenk’s result. In comparison, the resolution values in Table 3 rapidly and monotonically decrease for increasing  $n$ , so that *a fortiori* it may seem probable that they drop below 0.5 for some higher  $n$  (At resolution 0.5 each graph on average has a nondistinguished mate.).

It is tempting to ask similar questions with respect to the other graph invariant (combinations) considered here. The resolution of  $J$  for general graphs rapidly decreases for increasing  $n$  (Table 2), dropping to 0.31 for  $n = 10$  already, so that from this experimental point of view almost all graphs are  $J$ -equivalent (i.e. have a  $J$ -equivalent mate). The situation is less clear for  $J$  and the trees. Though  $J$  is still one of the best-discriminating simple invariants for trees (as we saw it is even better than the adjacency spectrum in this respect, Table 15), our data do not exclude the possibility that almost all trees are  $J$ -equivalent. This may seem paradoxical, but it is not a contradiction.

On the limited data obtained here for the distance spectrum and the  $J$ /spectrum combinations we do not want to speculate. Their resolutions also drop for increasing  $n$ , but slowly and not always monotonically, so that it seems possible but by no means clear that statements similar to the above are true for them.

#### ACKNOWLEDGMENT

Help with the handling of very large files and reformatting lists of connection tables by Dipl.-Math. J. Braun and Dipl.-Math R. Gugisch is gratefully acknowledged. We thank Dr. M. Grohe for ref 16.

#### REFERENCES AND NOTES

- (1) (a) Bonchev, D. The Overall Wiener Index – A New Tool for Characterization of Molecular Topology. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 582–592. (b) Bertz, S. H. Complexity of Molecules and their Synthesis. In *Complexity in Chemistry*; Bonchev, D., Rouvray, D. H., Eds.; Mathematical Chemistry Series, Vol. 7, in press. (c) Varmuza, K.; Scsibran, H. Substructure Isomorphism Matrix. *J. Chem. Inf. Comput. Sci.* **2000**, *40*, 308–313. (d) Poroikov, V. V.; Filimonov, D. A.; Borodina, Yu. V.; Lagunin, A. A.; Kos, A. Robustness of Biological Activity Spectra Predicting by Computer Program PASS for Noncongeneric Sets of Chemical Compounds. *J. Chem. Inf. Comput. Sci.* **2000**, *40*, 1349–1355.
- (2) (a) Blinov, K. A.; Elyashberg, M. E.; Molodtsov, S. G.; Williams, A. J.; Martirosian, E. R. An Expert System for Automated Structure Elucidation Utilizing  $^1\text{H}$ - $^1\text{H}$ ,  $^{13}\text{C}$ - $^1\text{H}$  and  $^{15}\text{N}$ - $^1\text{H}$  2D NMR Correlations. *Fresenius J. Anal. Chem.* **2001**, *369*, 709–714. (b) Varmuza, K.; Penchev, P. N.; Scsibran, H. Maximum Common Substructures of Organic Compounds Exhibiting Similar Infrared Spectra. *J. Chem. Inf. Comput. Sci.* **1998**, *38*, 420–427. (c) Seebass, B.; Pretsch, E. Automated Compatibility Tests of the Molecular Formulas or Structures of Organic Compounds with Their Mass Spectra. *J. Chem. Inf. Comput. Sci.* **1999**, *39*, 713–717.
- (3) Rücker, G.; Rücker, C. Automatic Enumeration of All Connected Subgraphs. *MATCH – Commun. Math. Comput. Chem.* **2000**, *41*, 145–149.
- (4) Rücker, G.; Rücker, C. On Finding Nonisomorphic Connected Subgraphs and Distinct Molecular Substructures. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 314–320; 865.
- (5) Rücker, G.; Rücker, C. Substructure, Subgraph, and Walk Counts as Measures of the Complexity of Graphs and Molecules. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 1457–1462.
- (6) Balaban, A. T. Highly Discriminating Distance-Based Topological Index. *Chem. Phys. Lett.* **1982**, *89*, 399–404.
- (7) (a) Balaban, A. T.; Quintas, L. V. The Smallest Graphs, Trees, and 4-Trees with Degenerate Topological Index  $J$ . *MATCH – Commun. Math. Comput. Chem.* **1983**, *14*, 213–233. (b) Razinger, M.; Chrétien, J. R.; Dubois, J. E. Structural Selectivity of Topological Indexes in Alkane Series. *J. Chem. Inf. Comput. Sci.* **1985**, *25*, 23–27.
- (8) (a) Bonchev, D.; Mekenyan, O.; Trinajstić, N. Isomer Discrimination by Topological Information Approach. *J. Comput. Chem.* **1981**, *2*, 127–148. (b) Balasubramanian, K.; Basak, S. C. Characterization of Isospectral Graphs Using Graph Invariants and Derived Orthogonal Parameters. *J. Chem. Inf. Comput. Sci.* **1998**, *38*, 367–373.
- (9) Randić, M.; Vracko, M.; Nović, M. Eigenvalues as Molecular Descriptors. In Diudea, M. V., Ed.; *QSPR/QSAR Studies by Molecular Descriptors*; Nova Science Publ.: Huntington, NY, 2001.
- (10) Balaban, A. T.; Harary, F. The Characteristic Polynomial Does Not Uniquely Determine the Topology of a Molecule. *J. Chem. Doc.* **1971**, *11*, 258–259.
- (11) Cvetković, D.; Rowlinson, P.; Simić, S. *Eigenspaces of Graphs*; Cambridge University Press: Cambridge, 1997.
- (12) Schwenk, A. J. Almost All Trees are Cospectral. In *New Directions in the Theory of Graphs*; Harary, F., Ed.; Academic Press: New York, 1973; pp 275–307.
- (13) Mihaljić, Z.; Veljan, D.; Amić, D.; Nikolić, S.; Plavšić, D.; Trinajstić, N. The Distance Matrix in Chemistry. *J. Math. Chem.* **1992**, *11*, 223–258.
- (14) Cvetković, D. M.; Doob, M.; Gutman, I.; Torgasev, A. *Recent Results in the Theory of Graph Spectra*; Elsevier: Amsterdam, 1988; p 128.
- (15) McKay, B. D. On the Spectral Characterisation of Trees. *Ars Combinatorica* **1977**, *3*, 219–232.
- (16) Cai, J.-Y.; Fürer, M.; Immerman, N. An Optimal Lower Bound on the Number of Variables for Graph Identification. *Combinatorica* **1992**, *4*, 389–410.
- (17) Shrikhande, S. S. On a Characterization of the Triangular Association Scheme. *Ann. Math. Statist.* **1959**, *30*, 39–47.
- (18) Weisfeiler, B. *On Construction and Identification of Graphs*; Lecture Notes in Mathematics No. 558, Springer: Berlin, 1976.
- (19) Mathon, R. *Proc. 9th S.-E. Conf. Combinatorics, Graph Theory and Computing* **1978**, 499.
- (20) Read, R. C.; Wilson, R. J. *An Atlas of Graphs*; Clarendon Press: Oxford: 1998; pp 4 and 7.
- (21) Recall that cubane, cuneane, and octabisvalene are pentacyclic octanes ( $n = 8$ ,  $m = 12$ ,  $c = 5$ ). A heptacyclic octane ( $n = 8$ ,  $m = 14$ ,  $c = 7$ ) would require two additional bonds in a cubane skeleton, for example.
- (22) Kerber, A.; Laue, R.; Grüner, T.; Meringer, M. MOLGEN 4.0. *MATCH – Commun. Math. Comput. Chem.* **1998**, *37*, 205–208.
- (23)  $J$  and eigenvalues were calculated as double precision numbers, for comparisons eight decimal places and seven decimal places were used for  $J$  and eigenvalues, respectively, for reasons detailed in ref 4.
- (24) We note in passing that having the necessary software we also determined the triplet of smallest isospectral 4-graphs ( $n = 8$ ,  $m = 10$ , corresponding to 2-methyltricyclo[3.2.0.0<sup>1,6</sup>]heptane, 3,4-dimethyltricyclo[3.1.0.0<sup>2,6</sup>]hexane, and [4.1.1]propellane) and the quadruplet of smallest isospectral graphs ( $n = 8$ ,  $m = 11$ , corresponding to 7-methyltetracyclo[2.2.1.0.1<sup>3,0</sup>1<sup>5</sup>]heptane, 6-methyltetracyclo[3.1.1.0.1<sup>3,0</sup>3<sup>5</sup>]heptane, tetracyclo[2.2.2.0.1<sup>3,0</sup>1<sup>4</sup>]octane, and 1-methyltetracyclo[3.2.0.0.1<sup>3,0</sup>2<sup>7</sup>]heptane). Recently an isospectral triplet of  $n = 9$ ,  $m = 16$  and an isospectral quadruplet of  $n = 9$ ,  $m = 19$  were published.<sup>26b</sup> Further we identified the pair of smallest distance-isospectral connected simple graphs, which have  $n = 7$ ,  $m = 10$ , corresponding to tetracyclo[3.1.1.0.1<sup>3,0</sup>3<sup>5</sup>]heptane and tetracyclo[2.2.1.0.1<sup>3,0</sup>1<sup>5</sup>]heptane.
- (25) In fact, this was not so surprising, after these graphs could not be pairwise distinguished using  $J$  and the distance eigenvalues.
- (26) (a) Hosoya, H.; Nagashima, U.; Hyugaji, S. Topological Twin Graphs. Smallest Pair of Isospectral Polyhedral Graphs with Eight vertices. *J. Chem. Inf. Comput. Sci.* **1994**, *34*, 428–431. (b) Hosoya, H.; Ohta, K.; Satomi, M. Topological Twin Graphs II. Isospectral Polyhedral Graphs With Nine and Ten vertices. *MATCH – Commun. Math. Comput. Chem.* **2001**, *44*, 183–200.
- (27) Despite having many graph invariants identical, such pairs of graphs can be differentiated by their chemical names, easily obtained using program POLCYC, e.g. **1** corresponds to heptacyclo[3.3.0.0.1<sup>3</sup>-.0.1<sup>4</sup>.0.2<sup>6</sup>.0.7<sup>0</sup>6<sup>8</sup>]octane, **2** to heptacyclo[4.2.0.0.1<sup>3</sup>.0.1<sup>7</sup>.0.2<sup>4</sup>.0.5<sup>0</sup>5<sup>8</sup>]octane.
- (28) If index values are compared in a sample of several or all edge numbers  $m$  within a constant vertex number  $n$ , i.e., without prior sorting by  $m$ , then additional degeneracies will occur. For instance, cyclooctane ( $n = 8$ ,  $m = 8$ ) and cubane ( $n = 8$ ,  $m = 12$ ) share the  $J$  value 2.000. Since NIMSG always sorts by  $n$  and  $m$ , we are not interested in such degeneracies here. It is well-known that graphs of different  $n$  can share the same  $J$  value, e.g. cyclohexane ( $n = 6$ ,  $m = 6$ ) also has  $J = 2.000$ .
- (29) Initially we were concerned that  $J$  and the distance spectrum, both being derived from the distance matrix, might tend to exhibit degeneracies for the same pairs of graphs. Fortunately, as foreseen already from the results shown in Figures 1 and 2, such concerns did not materialize to a large extent. However, the moderate improvement in the resolution of the distance spectrum on addition of  $J$  compared to the large improvement in the resolution of the adjacency spectrum on addition of  $J$  (Tables 15/16) may be interpreted to be partially due to such an effect.
- (30) The lower homologues of **13** and **14** lacking the methyl group are neither  $J$ -equivalent nor isospectral.
- (31) Grüner, T. Program GRADPART; Diploma Thesis, University of Bayreuth, 1995.
- (32) Balaban, A. T. Topological Index  $J$  for Heteroatom-Containing Molecules Taking into Account Periodicities of Element Properties. *MATCH – Commun. Math. Comput. Chem.* **1986**, *21*, 115–122.