

This is a postprint version of the following published document:

Sidorova, J., García, J. (2015). Bridging from syntactic to statistical methods: Classification with automatically segmented features from sequences. *Pattern Recognition*, 48, pp. 3749-3756.

DOI: [10.1016/j.patcog.2015.05.001](https://doi.org/10.1016/j.patcog.2015.05.001)

© Elsevier, 2015



This work is licensed under a [Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nc-nd/4.0/).

# Bridging from Syntactic to Statistical Methods: Classification with Automatically Segmented Features from Sequences

J. Sidorova<sup>a,b</sup>, J. Garcia<sup>a</sup>

<sup>a</sup> *Lab for Applied Artificial Intelligence,  
Faculty of Computer Science,  
Universidad Carlos III de Madrid.*

<sup>b</sup> *Department of Computer Science and Engineering  
Blekinge Institute of Technology,  
Sweden.  
julia.a.sidorova@gmail.com*

---

## Abstract

To integrate the benefits of statistical methods into syntactic pattern recognition, a *Bridging Approach* is proposed: (i) acquisition of a grammar per recognition class; (ii) comparison of the obtained grammars in order to find substructures of interest represented as sequences of terminal and/or non-terminal symbols and filling the feature vector with their counts; (iii) hierarchical feature selection and hierarchical classification, deducing and accounting for the domain taxonomy. The bridging approach has the benefits of syntactic methods: preserves structural relations and gives insights into the problem. Yet, it does not imply distance calculations and, thus, saves a non-trivial task-dependent design step. Instead it relies on statistical classification from many features. Our experiments concern a difficult problem of chemical toxicity prediction. The code and the data set are open-source.

*Keywords:* syntactic pattern recognition, grammatical inference, feature

*April 30, 2015*

## 1. Introduction

Statistical pattern recognition has a simple representation in the form of vectors allowing efficient ways to manipulate them, while syntactic pattern recognition has expressive representations, – graphs, strings, and so on, – but lacks object manipulation tools. Until recently, the syntactic and structural communities coexisted without much interaction. Yet, with the ever increasing difficulty of tasks in pattern recognition, more and more often the questions are asked: –*Can we have advantages of both paradigms?* –*Which are the trade-offs in such combinations?*

**Syntactic pattern recognition** can be used if there is a clear structure in the patterns and a *grammar* can be observed in a natural way. Forcing modeling on data, e.g. imposing linear ordering, hampers the performance [1]. Objects are represented by a variable-cardinality set of symbolic features.

Let there be  $n$  different grammars  $G_1, \dots, G_n$ , one for each recognition class  $C_k$   $k = 1, \dots, n$ . A *pattern*  $p_x$  of an object  $x$ , – where,  $x$  can be a written digit, speech sample, protein sequence, etc. – must first be transformed to a sequence of terminal symbols, that is, smallest units. For example, a protein sequence as a string

$$p_x = ATTTGGGGCTTATATAT, \quad (1)$$

where  $A, T, C, G$  are terminal symbols corresponding to the four nucleotides in the DNA. Examples of a recognition class  $C_k$  form a training set  $S(C_k)$ :

$$S(C_k) = \{p_{k_1}, p_{k_2}, p_{k_3}, \dots\}, \quad (2)$$

21 and a grammar  $G_k$  is sought, such that  $L(G_k) \supseteq S(C_k)$ . For a review of  
22 grammatical inference issues the reader is referred to [2], [3].

23 There exist various *distance metrics* to measure similarity between pat-  
24 terns. Let  $D(p_x, C_k)$  be some distance from a pattern  $p_x$  to a class  $C_k$ . The  
25 (smallest) distance between an input pattern  $p_x$  and a recognition class  $C_k$ <sup>1</sup>  
26 is

$$D(p_x, C_k) = \min\{D(p_x, p_k) | p_k \in L(G_k)\}. \quad (3)$$

27 In the literature, *three main approaches* to syntactic pattern recognition  
28 are typically singled out [4]:

- 29 – with an error-correcting parser,
- 30 – distance-based, and
- 31 – stochastic.

32 *An error-correcting parser* decides whether  $p_x$  belongs to  $L(G_i)$  or not.  
33 If  $p_x$  belongs to  $L(G_i)$ ,  $x$  is assigned to category  $C_i$ , and it is rejected other-  
34 wise. *The distance-based scheme* computes a distance from  $p_x$  to  $L(G_k)$ . If  
35  $D(p_x, L(G_i))$  is smallest among all the classes  $C_1 \dots C_n$ ,  $x$  is assigned to cate-  
36 gory  $C_i$ . Here, a statistical component is often added, and the distances to  
37 recognition classes are the input to a statistical classifier, where C4.5 or the  
38 kNN are known to perform well and keep the classification process human  
39 readable. *Stochastic schemes* consist in adding occurrence probabilities to  
40 productions in the schemes defined above.

41 Obviously, object representation is crucial, and *graphs* would be ideal in  
42 many applications, but learning graph grammars is largely infeasible due to

---

<sup>1</sup>or equivalently, between  $p_x$  and  $L(G_k)$

43 complexity issues<sup>2</sup>, instead graph embedding, e.g. [6], [7], and kernel meth-  
44 ods, e.g. [8], [9], are used. For the research trend on graphs in pattern  
45 recognition, the reader is referred to [10]. *Strings* are suitable, since a regu-  
46 lar or context-free grammar can be efficiently learnt and similarity measures  
47 calculated. If the target language is regular, hidden Markov models (HMMs)  
48 have been used in many applications [11]. For example, they are the main-  
49 stream tool to discover chromatin states [12], or protein regions [13] with  
50 distinct biological functions. The problem is that HMMs treat sequences  
51 as one-dimensional strings of independent, uncorrelated symbols. Although  
52 computationally convenient, this assumption is not structurally realistic [14],  
53 because many phenomena have more complex structure than regular: nat-  
54 ural language, palindrome structures in biology, and so on. Furthermore,  
55 once the target structure rises in terms of structural complexity from regular  
56 to context-free, one must make quite a number of task-dependent modeling  
57 decisions, and as a result applications become harder to design and reuse.  
58 Still, such efforts exist in optical character recognition [15], analysis of coro-  
59 nary artery images [16], in chemical biodegradability prediction [17], [18],  
60 and some other.

61 **Statistical pattern recognition** has a simple representation in the form  
62 of vectors and efficient ways to manipulate them [19]. It has gained a much  
63 greater popularity than the syntactic paradigm. Yet, faced with ever grow-  
64 ing difficulty of tasks, a recent tendency is to adapt ideas from syntactic

---

<sup>2</sup>A problem of parsing non-trivial graph languages is PSPACE-complete or NP-  
complete. Defining graph-grammars generating languages with a polynomial membership  
problem is an open problem [5].

65 methods. For example, in image understanding, ontologies are used for the  
66 loss function design: it is less of an error to take a cat for a dog, since both  
67 are animals, than a cat for a truck. In image tagging, structurally related  
68 features were shown to improve performance: if a ship has been detected,  
69 the probability for the the sea should be high. In graph matching, structural  
70 information allows for constraint formulation: if a face is adjacent to a neck  
71 in one graph, it should be so in the other one, too. For an overview the  
72 reader is referred to [20]. Another idea proposed is to gain interpretability of  
73 predictive models in some creative task-dependent way, which often comes  
74 with a cost in recognition accuracy compared to black-box solutions or may  
75 require that the underlying linear model works well on the data set: for ex-  
76 ample, adding a heat map coloring technique to interpret linear support  
77 vector machine models [21].

78 This work, too, explores connections between the two paradigms, but our  
79 idea is different. In our previous work [18], we departed from the fact that  
80 there is a grammar for chemicals, very much like a natural grammar, and, we  
81 designed a syntactic pattern recognition scheme together with a procedure  
82 to search for important substructures in the grammars. In this submission,  
83 we propose to fill the feature vector with the counts of potentially important  
84 substructures. These substructures are automatically segmented, have an  
85 automatically chosen degree of structural abstraction and special statistical  
86 properties. The proposed *Bridging Approach* brings the following benefits:

- 87 1. The method’s essential capacity is *to cope in the absence of expert*  
88 *knowledge*, that is, no indications with respect to which features to  
89 extract or where to look for them in the input sequence.

- 90 2. *It gives insights into the problem* in two respects. Firstly, the method  
91 works with a variable-length parsable input and finds the regions of  
92 interest in sequences with a suitable level of abstraction for their rep-  
93 resentation. Secondly, subsequent hierarchical vector-based feature se-  
94 lection and classification account for the domain’s taxonomy.
- 95 3. *It is easier-to-implement* than a classical syntactic scheme, since it does  
96 not imply distance calculations. Therefore, it saves a non-trivial design  
97 step from the syntactic paradigm.

98 Our experiments concern a difficult problem of chemical toxicity prediction.  
99 Our parser processes molecules in the SMILES format, which is a string  
100 representation of a 2D molecular graph. From two sets of molecules with  
101 opposite properties  $S(G_{\oplus})$  and  $S(G_{\ominus})$ , a predictive model is built with *the*  
102 *Bridging Approach*.

103 The rest of the paper is organized as follows. Section 2 explains how  
104 chemicals are represented as strings and how they are parsed. Section 3  
105 explains the steps of *the Bridging Approach*. Section 4 covers the exper-  
106 iment. Finally, conclusions are drawn in Section 5. The SMILES parser  
107 and *the bridging approach* are available on request from the correspond-  
108 ing author. The database used for experiments is NCTRER DSSTOX at  
109 [http://www.epa.gov/nheerl/dsstox/sdf\\_nctrer.html](http://www.epa.gov/nheerl/dsstox/sdf_nctrer.html)

## 110 2. Parsing Chemicals

111 **The chemical language SMILES** was designed “*to represent molecular*  
112 *structure by a linear string of symbols, similar to a natural language*“ [22].  
113 A sequence in SMILES represents a molecular structure as a graph.

114 *Atoms:* Atoms are represented by their atomic symbols: C, Cl, N, O, etc.  
115 This is the only required use of letters in SMILES. Hydrogen atoms (H) are  
116 normally omitted, since valences make it clear where they are missing. For  
117 example, an atomic chain CCSCCCCC<sup>3</sup> is depicted in Figure 1.

118 *Bonds:* Single bonds are usually omitted in SMILES. Double and triple  
119 bonds are represented by the symbols = and #, respectively, for example, in  
120 Figure 2.

121 *Branches.* Branches are specified by enclosures in parentheses, as in Fig-  
122 ure 3.

123 *Cyclic Structures:* Cyclic structures are represented by breaking one sin-  
124 gle (or aromatic) bond in each ring. The bonds are numbered in any order,  
125 designating ring opening (or ring-closure) bonds by a digit immediately fol-  
126 lowing the atomic symbol at each ring closure. This leaves a connected  
127 noncyclic graph, which is written as a noncyclic structure, as in Figure 4.

128 With the rules above almost all organic structures can be described as  
129 strings. For more details, the reader is referred to [22].

130 **A context-free parser based on the SMILES grammar** we devel-  
131 oped creates a syntax tree from SMILES, see Appendix A for further details.

### 132 **3. The Bridging Approach**

133 Input is parsed structured data, the *Bridging Approach* will study it and  
134 build a predictive model based on its conclusions. Briefly, its steps are:

- 135 1. acquisition of a grammar per recognition class;

---

<sup>3</sup>Due to chemical convention in graphics, whenever a label on graph node is missing, it is C and a line segment represents a chemical bond.



- 136 2. comparison of the obtained grammars in order to find substructures  
137 of interest represented as sequences of terminal and/or non-terminal  
138 symbols and filling the feature vector with their counts;
- 139 3. hierarchical feature selection and hierarchical classification, deducing  
140 and accounting for the domain taxonomy.

141 **Step 1: acquisition of a grammar per recognition class.**

142 The general assumption is that objects with similar structures have simi-  
143 lar properties. Given two sets of examples from opposite classes (for example,  
144 active and non-active chemicals), we can learn grammars that account for  
145 their structures:  $L(G_k) \supseteq S(C_k)$ . Examples are taken from the training set  
146 one by one. Whenever an example cannot be parsed with the current gram-  
147 mar, the grammar is extended with new rules to accommodate the example.

148 The grammar inference algorithm from SMILES [18] is reproduced in Ap-  
149 pendix B. The input to it is a training set with parsed SMILES of chemicals  
150 belonging to the same activity class, and the output is the grammar  $G$  and  
151 table  $T$  of two columns:

152  $\langle$ production  $p$  $\rangle$  and  $\langle$ how many times  $p$  was used $\rangle$ .

153 **Step 2: comparison of the obtained grammars in order to re-**  
154 **veal substructures of interest represented as sequences of terminal**  
155 **and/or non-terminal symbols.**

156 For a binary problem<sup>4</sup>, the tables for the two classes,  $T_{\oplus}$  and  $T_{\ominus}$ , are  
157 compared, in order to search for the substructures of interest. There is a

---

<sup>4</sup>A multiclass problem can be recast into a series of binary classification problems with one-versus-all [23], one-versus-one [24] and error-correcting output codes [25], [26].

158 qualitative and quantitative aspect to this search. The qualitative aspect  
159 concerns feature segmentation: what are the substructures of interest and  
160 how they are represented as terminal and/or non-terminal symbols. The  
161 quantitative question is which statistical properties the substructures should  
162 possess, in order that the counts of their occurrences can be useful as features  
163 in statistical classification.

164 The qualitative issue is resolved by the grammar. Consider examples of  
165 productions:

$$166 \quad sig_2 \rightarrow sig_6 sig_6, \quad (4)$$

$$167 \quad sig_6 \rightarrow C1Csig_3CCC1, \quad (5)$$

$$168 \quad sig_6 \rightarrow C1CCCCC1. \quad (6)$$

168 The left-hand side of productions entirely depends on the right hand side  
169 and is redundant, that is, it is of the form  $sig_{arity}$ , where the arity is the  
170 number of units that appear in the right-hand side. Thus, we can work with  
171 the right-hand side only. The grammar defines how the substructures are  
172 segmented and the level of abstraction. In our example the substructures are  
173  $sig_6 sig_6$ ,  $C1Csig_3CCC1$ ,  $C1CCCCC1$ .

174 From a quantitative perspective, naturally one would look closely at fre-  
175 quently encountered molecular substructures that are *exclusive* for one class.  
176 Unfortunately, such ideal "structural alerts" are infrequent due to many  
177 chemical exceptions, and we can't hope that they alone can solve the classifi-  
178 cation and explanatory tasks. *Common substructures* need to be considered.  
179 In order to favor the ones that are more frequent in one class and less frequent  
180 in the other, the ones that have the importance value greater than average

181 are taken, where

$$Importance(X_i) = |count(X_i) \text{ in } T_{\ominus} - count(X_i) \text{ in } T_{\oplus}|. \quad (7)$$

182 **Step 3: hierarchical feature selection from the pool of substructures**  
183 **of interest and hierarchical classification, deducing and ac-**  
184 **counting for the domain taxonomy at the feature selection and**  
185 **classification steps.**

186 Initially, two types of substructures are filtered: the substructures found  
187 in one of the classes exclusively and common substructures that are more  
188 frequent to in class than in the other. In order to incorporate this intuition  
189 into a predictive model, it needs to be backed with statistics. Additionally:

- 190 1. It should be taken into account that many domains have natural tax-  
191 onomies, for example species, chemicals etc form families and subfami-  
192 lies. Within a taxonomic category, objects have comparable structures  
193 and property-specific structural clues can further be discovered. In  
194 terms of structure (morphology), the gold fish can be compared to the  
195 carp, but not to the hamster.
- 196 2. The method should keep the criteria for classification human-readable.

197 Decision trees are a standard choice, when human readability and gaining  
198 insights are sought. Further, C4.5 [27] automatically partitions the feature  
199 space and chooses appropriate features for classification in each subregion.

## 200 4. Experiment

201 **Data:** A large number of chemicals present in the environment are es-  
202 trogens, that is they are structurally similar to hormones and disrupt en-

203 docrine functions in animals and humans [28]. The NCTR (National Cen-  
204 ter for Toxicological Research) Estrogen Receptor Binding database [29],  
205 [http://www.epa.gov/nheerl/dsstox/sdf\\_nctrer.html](http://www.epa.gov/nheerl/dsstox/sdf_nctrer.html), consists of 232 chemi-  
206 cals. Its creation was motivated by the desire to summarize the knowledge  
207 about estrogens and have a reliable data set of consistent design that would  
208 fully cover structurally diverse set of natural, synthetic, and environmental  
209 estrogens. Once the list of chemicals had been composed by experts, they  
210 were tested on rats in well validated and standardized analytical procedure.  
211 The estrogen activity was measured on the scale from 0 to 100: 0 corresponds  
212 to *inactive*, the chemicals with the activity values  $\geq 23$  are labeled as *active*,  
213 and the structures labeled as *inconclusive* have the activity value equal to  
214 5. The authors also provided a set of chemical rules linking substructures,  
215 types of chemicals and their resulting activity. Among the 232 samples: 89  
216 chemicals are active (the  $\oplus$  class), and 123 are inactive (the  $\ominus$  class), and 8  
217 chemicals are labeled as *inconclusive*. We decided to include the inconclusive  
218 chemicals as a third class to observe the tendencies.

219 **Learning settings:** the experiments were carried out in 10-fold cross  
220 validation. A 90% part of data was taken for training purposes to carry out  
221 steps 1-3 of the *Bridging Approach*. The remaining 10% was used for testing.

222 **Results:** The confusion matrix for the experiment is presented as Table  
223 1. As had been expected, the grammar for the inconclusive class was very  
224 small (since it had too few training examples) and therefore useless. Conse-  
225 quentially, the inconclusive class could not be recognized, and its chemicals  
226 appeared to fall randomly into the active and inactive classes. Further, when  
227 calculating recognition accuracy and other characteristics, the inconclusive

228 class was not taken into account. The overall recognition is 75%, recall =  
229 0.69, precision = 0.71 and F-measure = 0.7.

230 A predictive model is considered successful, if its accuracy is better than  
231 70% [30]. Our result is above this baseline, is comparable to some studies on  
232 the same database (67% [31], 78% [32], 79% [33]), but is considerably behind  
233 the best result reported of 85% [34] with two black-box methods that imple-  
234 mented the Random Forest with 500 trees and Classification by Ensembles  
235 from Random Partitions (CERP). That is in line with the literature, e.g.  
236 [34]: readability often goes with a cost in accuracy.

237 Once the recognition capacity of the method had been concluded to be  
238 satisfactory, the learning procedure was repeated on the whole dataset to  
239 obtain a decision tree that summarizes the activity in terms of structural  
240 features, which is depicted in Figure 12 in Appendix C.

241 The database creators provided an expert model with *if-then*-rules to  
242 summarize extrogenic activity based on advanced expertise in organic chem-  
243 istry [29], for details see Appendix C. We compared the method’s findings  
244 against the expert model. The expert rules were made following all the con-  
245 ventions and accommodating systematic chemical theory. The data-driven  
246 model had a limited data set with structures and labels only. Yet by far and  
247 large, the data-driven conclusions are in line with the expert rules. Some-  
248 times, the expert model uses parameters, other than presence/absence of a  
249 substructure, for example, a solubility-related coefficient  $\log p$ . These pa-  
250 rameters, too, can be successfully predicted from structure, and the *Bridging*  
251 *Approach* copes in the absence of this knowledge. Further details on the  
252 model comparison are given in Appendix C.

## 253 5. Conclusions

254 Our idea has been to bridge from syntactic to statistical pattern recog-  
255 nition. The feature vector is filled with the counts of the substructures that  
256 are extracted from grammars. Having grammars automatically solves the  
257 task of feature segmentation and the choice of degree of abstraction for their  
258 representation. The selected features have sophisticated statistical proper-  
259 ties, that is, max information gain at a particular point of the hierarchically  
260 divided feature space.

261 Compared to the syntactic paradigm, the new traits are:

- 262 • The proposed bridging model is directly recyclable in other applica-  
263 tions, as long as the input can be parsed.
- 264 • It does not imply distance calculations and, instead, relies on vector  
265 classification, and, therefore, saves a non-trivial design step compared  
266 to the syntactic paradigm.

267 Having gained the new advantages, the method preserves the inherent strengths  
268 of the syntactic paradigm:

- 269 • Its essential capacity is to cope in the absence of expert knowledge,  
270 that is, no indications which features to extract.
- 271 • It preserves structural relations and works with a variable-length parsable  
272 input. It finds regions of interest in sequences with a suitable level of  
273 abstraction for their representation, and learns a decision tree that op-  
274 erates on presence/absence of these structures. Altogether, it leads to  
275 human-readable classification and gives insights into the problem.

276 **6. Conflict of Interests**

277 None declared.

278 **7. Acknowledgements**

279 As usually, we thank Torben Hagerup and Ricard Gavalda for thoughtful  
280 advice and guidance, – part of the work was completed while JS visited the  
281 LARCA at UPC. We acknowledge useful discussions with Antonio Berlanga,  
282 Florian Leitner, and Andrew Moss. JS acknowledges the *estancias postdoc-*  
283 *torales* grant at the UC3M.

284 **Appendix A. Parsing examples**

285 Input to the parser is a SMILES of a chemical compound. The parser  
286 starts with the first atom in the SMILES string, uniquely identifies each  
287 atom with its position number and disambiguates which atom is linked to  
288 which other atoms and by which type of bond. Then, it reconstructs a tree  
289 representation of the compound.

290 An obvious challenge is that different SMILES exist for the same molecule.  
291 For example, a molecule from Figure 1 can be rotated and written as CCC-  
292 CCSCC in place of CCSCCCCC. Canonical SMILES are not a solution, since  
293 they can't be drawn for the reason that are a hash value due to principles of  
294 their construction, and we don't want a black box construction. Our solu-  
295 tion is sorting substructures in a natural order, when comparing sequences  
296 of substructures in grammar inference and in search.

297 Generally, the parser implements the SMILES language. Given the de-  
298 scription from Section 2, few additional decisions are left to be made with

299 respect to how non-terminals are assembled and finally reduced to the start  
300 symbol  $S$ . The additional rewriting rules are as follows.

301 *Rule 1:* under our modeling non-terminals are denoted with the symbol  
302  $sig_{arity}$ <sup>5</sup> and differ with respect to their *arity*, which is the number of the  
303 non-terminal’s child nodes, for example:

304



305 The numbers do not count, since they are special symbols.

306 *Rule 2:* Atomic chains, that is, molecules without rings or branches, are  
307 reduced directly to the start symbol  $S$ :



308 Figure 1 depicts this molecule<sup>6</sup>.

309 *Rule 3:* unlike atomic chains, branched and cycles are reduced to a non-  
310 terminal. For example, an atom and a branch hanging from it is reduced  
311 to a non-terminal of a corresponding arity. CC(CCCBr)CC is reduced to  
312  $Csig_5CC$ , where  $sig_5 \rightarrow C(CCCBr)$ , as in Figure 5.

313 An example of multiple branches stemming from the same atom is CC(F)(Br)I,  
314 depicted in Figure 6.

315 *Rule 4:* children nodes of a non-terminal node can be atoms and/or  
316 substructures, and in order to calculate arity the number of such units is  
317 counted.

---

<sup>5</sup>Traditionally non-terminals are labeled with  $\sigma_{arity}$ , which we spell in the Latin alphabet as  $sig_{arity}$ .

<sup>6</sup>Also a functionality to depict SMILES can be useful, e.g. [35]



318 *Rule 5:* if a ring is not a stand-alone ring as in Figure 2, the starting  
319 atom of the cycle (the one after which the number is put) is marked as *sig*<sub>1</sub>,  
320 and this disambiguates the branch from a cycle.

321 Very similar SMILES can lead to different parsing results, an example  
322 of the significance of parentheses is CC(C1CCC1)CC in Figure 7 (left) and  
323 CCC1CCC1CC in Figure 8.

324 *Rule 6:* In parser’s output, shared atoms have a special tag that they  
325 belong to two different cycles. Otherwise, for example C2OC1CCC2CC1 in  
326 Fig 9, is simply reduced to *sig*<sub>6</sub>*sig*<sub>6</sub>.

327 The above examples were simple to illustrate the parsing decisions. A  
328 couple of more complex molecules from the DSSTox NCTRER database are  
329 drawn in Figures 9 and 10.

## 330 Appendix B. Grammar Inference Algorithm

331 Under our modeling non-terminals are all marked with the symbol *sig*  
332 and differ with respect to their *arity*, which is the number of their child  
333 nodes, for example:

334

$$sig_6 \rightarrow N1CCCCC1. \quad (B.1)$$

335 The parsed SMILES are processed in postorder. In the algorithm below:

336 – *j* is the number of a node in post-order enumeration;

337 – *X*<sub>*j*</sub> is a string of the child nodes of the node *j*:

$$X_j = x_1x_2\dots x_l, \quad (B.2)$$

338 with  $l \geq 1$ . For example above, for *sig*<sub>6</sub> the string of child nodes is *NCCCCC*.

339 Since the tree graph is traversed in post-order, at the point of reducing *X*<sub>*j*</sub> to

340 a non-terminal  $sig_l$ , each of its child nodes  $x_1x_2\dots x_l$  have been parsed either  
341 as atoms or as non-terminal  $sig_{arity}$  nodes.

**Data:** The training set  $D$  of size  $n$  with parsed SMILES of chemicals belonging to the same activity class.

**Initialization:** Set  $G$  to contain empty sets for

- the set of atoms  $A$ ,
- the set of non-terminals  $N$ ,
- the set of rules  $P$ ,
- the start symbol  $S$ ,

and an empty table  $T$  with two columns:  $\langle p: \text{production from } P \rangle \langle \text{count}(p): \text{how many times } p \text{ was used} \rangle$ .

$n =$  the number of instances in  $D$ .

**for**  $i = 0$  **to**  $n$ , **while**  $i < n$  **do**

$i = i + 1$ ;

    read the  $i^{\text{th}}$  SMILES in  $D$ ;

    in postorder, **for** each node  $j$  in SMILES **do**

**If** any atoms from  $X_j$  are not in  $A$ , add them to  $A$ .

**If** the string  $X_j$  can not be reduced with productions from  $P$

        add the rule:  $\text{sig}_l \mapsto x_1x_2\dots x_l$  to  $P$ ;

**If**  $\text{sig}_l$  is not in  $N$ , add it to  $N$ .

        }

        Let  $p$  be the rule used to reduce  $X_j$ ;

**if**  $p$  is not in  $T$ , add  $p$  to  $T$  with  $\text{count}(p) = 0$ .

$\text{count}(p) = \text{count}(p) + 1$ ;

**end**

**end**

**Result:** the grammar  $G = (A, N_1S, P)$  generalizing the activity class to which the input samples in  $D$  belong and a table  $T$  with grammar productions and their counts.

**Algorithm 1:** Polynomial time algorithm for grammatical inference of structures belonging to the same activity class from their parsed SMILES.

## 343 Appendix C. Data-driven and Expert Model

344 **The data-driven** model in the form of *if-, then-* rules obtained with  
345 the *Bridging Approach* is depicted in Figure 12. The conditions check for  
346 the presence of particular structural alerts. The ratio at the leaf boxes is  
347  $\langle$ number of correctly classified samples, number of errors $\rangle$ .

348 **The expert model** from [29] is summarized to:

- 349 1. If a chemical contains no ring structure, it is unlikely to be an estrogen  
350 receptor ligand (ER-ligand).
- 351 2. If a chemical has a nonaromatic ring structure, then it is unlikely to be  
352 an ER ligand, if it does not contain an O, S, N.
- 353 3. If a chemical has a non-OH aromatic structure, then its binding po-  
354 tential is dependent on the existence of key structural features and a  
355 solubility-related coefficient  $\log p$ .
- 356 4. If a chemical contains a phenolic ring, then it tends to be an ER ligand,  
357 if it contains any additional key structural features. For the chemicals  
358 containing a phenolic ring separated from another benzene ring with  
359 the number of bridge atoms ranging from none to three, it will most  
360 likely be an ER ligand.

361 The rules in the expert model are based on systematic chemical knowl-  
362 edge. The data-driven model had a limited data set with chemical structures  
363 and labels for their activity only. Yet, by and large, the data-driven conclu-  
364 sions are in line with the expert rules:

- 365 1. Many of the chemicals covered by the 1st rule of the expert model  
366 end up at node 26 and node 31 passing as negative through numerous  
367 check-ups on the presence of different cyclic substructures.

- 368 2. The chemicals that satisfy the expert rule 7 from the original diagram  
369 [29] end up at the nodes 16 and 11. Node 9 ( =C@, O) is equivalent to  
370 the presence of a phenolic ring.
- 371 3. The expert model has complex cases where the binding potential is  
372 determined with the help of non-structural information such as *log p*.  
373 The data-driven model is not allowed to use any additional information  
374 and accommodates these chemicals checking a lengthy list of structural  
375 conditions.

## 376 References

- 377 [1] M. Venguerov, P. Cunningham, *Generalized syntactic pattern recognition*  
378 *as a unifying approach in image analysis*, in: The 2nd international work-  
379 shop on statistical techniques in pattern recognition, 1998, pp 913-920.
- 380 [2] R.G. Parekh, V. Honavar, *Grammar inference, automata induction, and*  
381 *language acquisition*, in: R. Dale, H. Moisl, H. Somers (Eds.), Handbook  
382 of Natural Language Processing, 2000, pp 727-764.
- 383 [3] C. De la Higuera, *A bibliographical study of grammatical inference*, Pat-  
384 tern Recognition 38 (2005) 1332-1348.
- 385 [4] E. Tanaka, *Theoretical aspects of syntactic pattern recognition*, Pattern  
386 Recognition 28 (1995) 1053-1061.
- 387 [5] M. Flasiński, J. Jurek, *Fundamental methodological issues of syntactic*  
388 *pattern recognition*, Pattern Anal. Appl. 17 (2014) 465-480.

- 389 [6] M. Ferrer, E. Valveny, F. Serratoso, H. Bunke, *Generalized median graph*  
390 *computation by means of graph embedding in vector spaces*, Pattern Recog-  
391 nition 43 (2010) 1642-1655.
- 392 [7] B. Luo, R.C. Wilson, E.R. Hancock, *Spectral embedding on graphs*, Pat-  
393 tern Recognition 36 (2003) 2213-2230.
- 394 [8] H. Bunke, K. Riesen, *Recent advances in graph-based pattern recogni-*  
395 *tion with applications in document analysis*, Pattern Recognition 44 (2011)  
396 1057-1067.
- 397 [9] B. Gauzere, P.A. Grenier, L. Brun, D. Villemin, *Treelet kernel incorpo-*  
398 *rating cyclic, stereo and inter pattern information in chemoinformatics*,  
399 Pattern Recognition 48 (2015) 356-367.
- 400 [10] M. Vento, *A long trip in the charming world of graphs for Pattern Recog-*  
401 *niton*, Pattern Recognition 48 (2015) 291-301.
- 402 [11] J. Garcia, R. Aler, J. Sidorova, *Data Analysis*, Open  
403 Courseware Universidad Carlos III de Madrid (2014)  
404 <https://www.youtube.com/watch?v=q4IHnJws3qE> retrieved on  
405 27/01/2015.
- 406 [12] J. Ernst, M. Kellis, *Discovery and characterization of chromatin states*  
407 *for systematic annotation of the human genome*, Nat. Biotechnol. 28 (2010)  
408 817-825.
- 409 [13] A. Arcas, I. Cases, A.M. Rojas, *Serine/threonine kinases and E2-*  
410 *ubiquitin conjugating enzymes in Planctomycetes: unexpected findings*, An-  
411 tonie Van Leeuwenhoek 104 (2013) 509-20.

- 412 [14] R. Durbin, S.R. Eddy, A. Krogh, G. Mitchison, *Biological sequence anal-*  
413 *ysis. Probabilistic models of proteins and nucleic acids*, Cambridge, UK,  
414 1999.
- 415 [15] J.M. Sempere, D. Lopez, *Learning Decision Trees and Tree Automata for*  
416 *a Syntactic Pattern Recognition Task*, in: F.J. Perales et al. (Eds.), *Pattern*  
417 *Recognition and Image Analysis*, Springer-Verlag Berlin Heidelberg 2003,  
418 pp. 943-950.
- 419 [16] M.R. Ogiela, R. Tadeusiewicz, *Syntactic reasoning and pattern recogni-*  
420 *tion for analysis of coronary artery images*, *Artif. Intell. Med.* 26 (2002)  
421 145-159.
- 422 [17] J. Sidorova, A. Fernandez, J. Cester, R. Rallo, F. Giralt, *Predicting*  
423 *biodegradable quality of chemicals with the TGI+. 3 classifier*, in: *The*  
424 *11th IASTED International conference on artificial intelligence and appli-*  
425 *cations*, 2011, pp 108-115.
- 426 [18] J. Sidorova, M. Anisimova, *NLP-inspired structural pattern recognition*  
427 *in chemical application*, *Pattern Recognit. Lett.* 45 (2014) 11-16.
- 428 [19] I.H. Witten, E. Frank, M.A. Hall, *Data Mining: Practical Machine*  
429 *Learning Tools and Techniques*, 3rd edition, Morgan Kaufmann Publish-  
430 ers, 2011.
- 431 [20] T. Caetano, *The interplay of statistical and structural pattern recog-*  
432 *niton from a machine learning perspective*, in: *International Confer-*  
433 *ence on Pattern Recognition Applications and Methods (ICPRAM)*, 2012,  
434 <http://vimeo.com/38450616>, retrieved on 27/01/2015.

- 435 [21] L. Rosenbaum, G. Hinselmann, A. Jahn, A. Zell, *Interpreting linear sup-*  
436 *port vector machine with heat map molecule coloring.* Journal of Chemoin-  
437 formatics 3 (2011): 11.
- 438 [22] D. Weininger, *SMILES, a chemical language and information system.*  
439 *1. Introduction to methodology and encoding rules,* J. Chem. Inf. Comput.  
440 Sci. 28 (1988) 31-36.
- 441 [23] R. Anand, K. Mehrotra, C.K. Mohan, S Ranka, *Efficient classification*  
442 *for multiclass problems using modular neural networks,* IEEE Transactions  
443 on Neural Networks 6 (1995) 117-124.
- 444 [24] T. Hastie, R. Tibshirani, *Classification by pairwise coupling,* Ann.  
445 Statist. 26 (1998), 451-471.
- 446 [25] T.G. Dietterich, G. Bakiri, *Solving multiclass learning problems via*  
447 *error-correcting output codes,* J. Artif. Intell. Res. 2 (1995) 263-286.
- 448 [26] M.A. Bagheri, Q. Gao, S. Escalera, *A genetic-based subspace analysis*  
449 *method for improving error-correcting output coding,* Pattern Recognition  
450 46 (2013) 2830-2839.
- 451 [27] R. Kovahi, J.R. Quinlan, *Data mining tasks and methods: Classifica-*  
452 *tion:Decision tree discovery,* in: W. Klsgen, J.M. Zytkow (Eds.), Hand-  
453 book of data-mining and knowledge discovery, Oxford University Press,  
454 Inc., New York, 2002, pp. 267-276.
- 455 [28] B. Hileman, *Hormone disrupter research expands,* Chem. Eng. News 75  
456 (1997) 24-25.



- 457 [29] H. Fang, W. Tong, L.M. Shi, R. Blair, R. Perkins, W. Branham, B.S.  
458 Hass, Q. Xie, S.L. Dial, C.L. Moland, D.M. Sheehan, *Structure-activity re-*  
459 *lationships for a large diverse set of natural, synthetic, and environmental*  
460 *estrogens*, Chem. Res. Toxicol. 14 (2001) 280-294.
- 461 [30] R. Hannu, *In silico toxicology – non-testing methods*, Front. Pharmacol.  
462 2 (2011) 33.
- 463 [31] N. Fonseca, V. Santos Costa, R. Camacho, *k-RNN: k-Relational nearest*  
464 *neighbor algorithm*, in: Symposium on applied computing (SAC), ACM,  
465 2008, pp. 944-948.
- 466 [32] N. Landwehr, A. Passerini, L. De Raedt, P. Frasconi, *k-FOIL: learn-*  
467 *ing relational kernels*, in: The 16th International Conference on Inductive  
468 Logic Programming (IPL), Short Papers (Muggleton, S. and Otero, R.,  
469 eds.), 2006, pp. 125-127.
- 470 [33] A. Karwath, L. De Raedt, *SMIREP: predicting chemical activity from*  
471 *SMILES*, J. Chem. Inf. Model. 46 (2006) 2432-2444.
- 472 [34] H. Ahn, H. Moon, M.J. Fazzari, N. Lim, J.J. Chen, R. Kodell, *Classi-*  
473 *fication from ensembles from random partitions of high-dimensional data*,  
474 Comput. Stat. Data An. 51 (2007) 6166-6179.
- 475 [35] N.M. OLBoyle, M. Banck, C.A. James, C. Morley, T. Vandermeersch,  
476 G.R. Hutchison, *Open Babel: An open chemical toolbox*, J Cheminf, 3  
477 (2011) 33.

**Julia Sidorova** is a postdoctoral fellow at the Universidad Carlos III de Madrid, graduated in 2002 from Moscow State University and got her PhD from Universitat Pompeu Fabra in 2009. Research interests include pattern recognition, chemical activity prediction and affect recognition from voice.

**Jesus Garcia** is an associate professor at the Universidad Carlos III de Madrid, got his PhD from Universidad Complutense de Madrid. Research interests include Data and Information fusion, Multiagent systems and distributed sensor networks, Machine learning and data mining, Applied computational intelligence in engineering: Air Traffic Control, Surveillance and Machine Vision, Autonomous Vehicles Navigation & Control.

**Table 1**

Predicted as --	Predicted as inconclusive	Predicted as +	Real Class:
66	0	27	--
5	0	3	inconclusive
29	0	100	+

Table 1: Confusion matrix: active (+), inconclusive, and inactive (-).

Figure 1



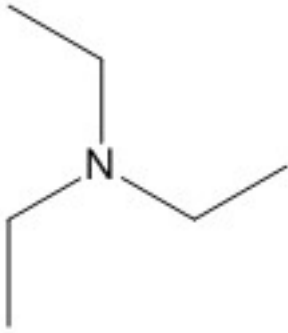
Atomic chain CCSCCCCC.

Figure 2



Double bond: C=C.

Figure 3



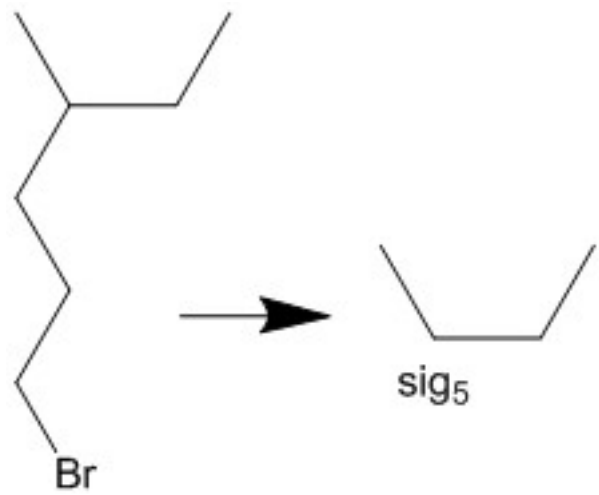
Branches: CCN(CC)CC.

Figure 4



SMILES: O1CCCCC1N1CCCCC1. Rings are broken, and a number is put to leave a mark where the bond was broken.

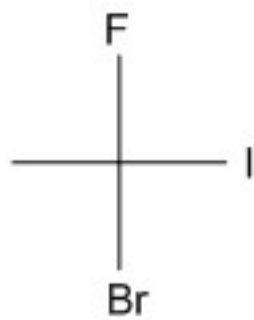
Figure 5



Branch is reduced to a non-terminal.

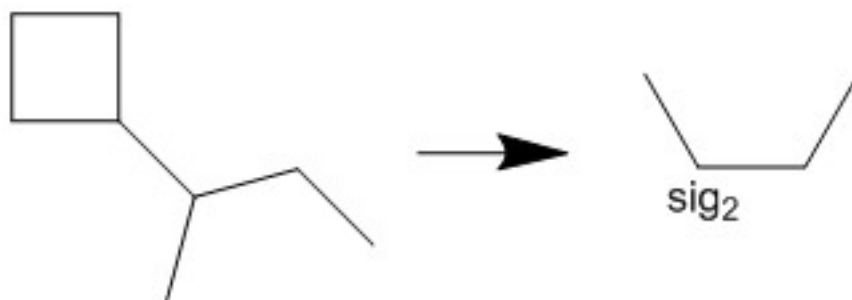


Figure 6



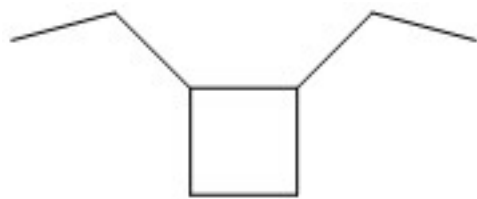
CC(F)(Br)I is reduced to "Csig3I".

Figure 7



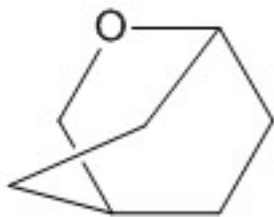
SMILES CC(C1CCC1)CC. The cycle is reduced to a non-terminal `sig2`. The nonterminal has two child units: an atom and another non-terminal.

Figure 8



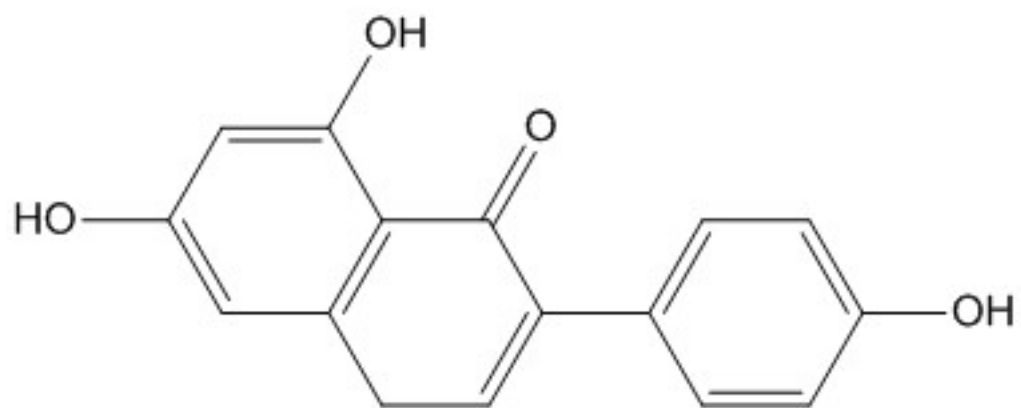
SMILES: CCC1CCC1CC.

Figure 9



Two intersecting rings: C2OC1CCC2CC1.  
In the parser's internal presentation the atoms there are special tags @1 (and @2) after each atom, disambiguating to which cycle it belongs. Shared atoms are followed by @1@2:  
C<sub>1</sub>@1@2O<sub>2</sub>@1@2C<sub>3</sub>@1@2C<sub>4</sub>@2C<sub>5</sub>@2C<sub>6</sub>@1@2C<sub>7</sub>@1C<sub>8</sub>@1. Subscript is atom's ID from the SMILES strings.

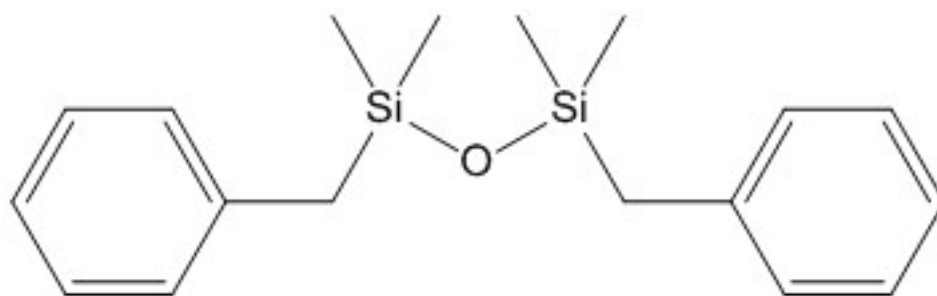
Figure 10



A molecule from the database:

O=C(C(C(C=C3)=CC=C3O)=CO2)C1=C2C=C(O)C=C1O.

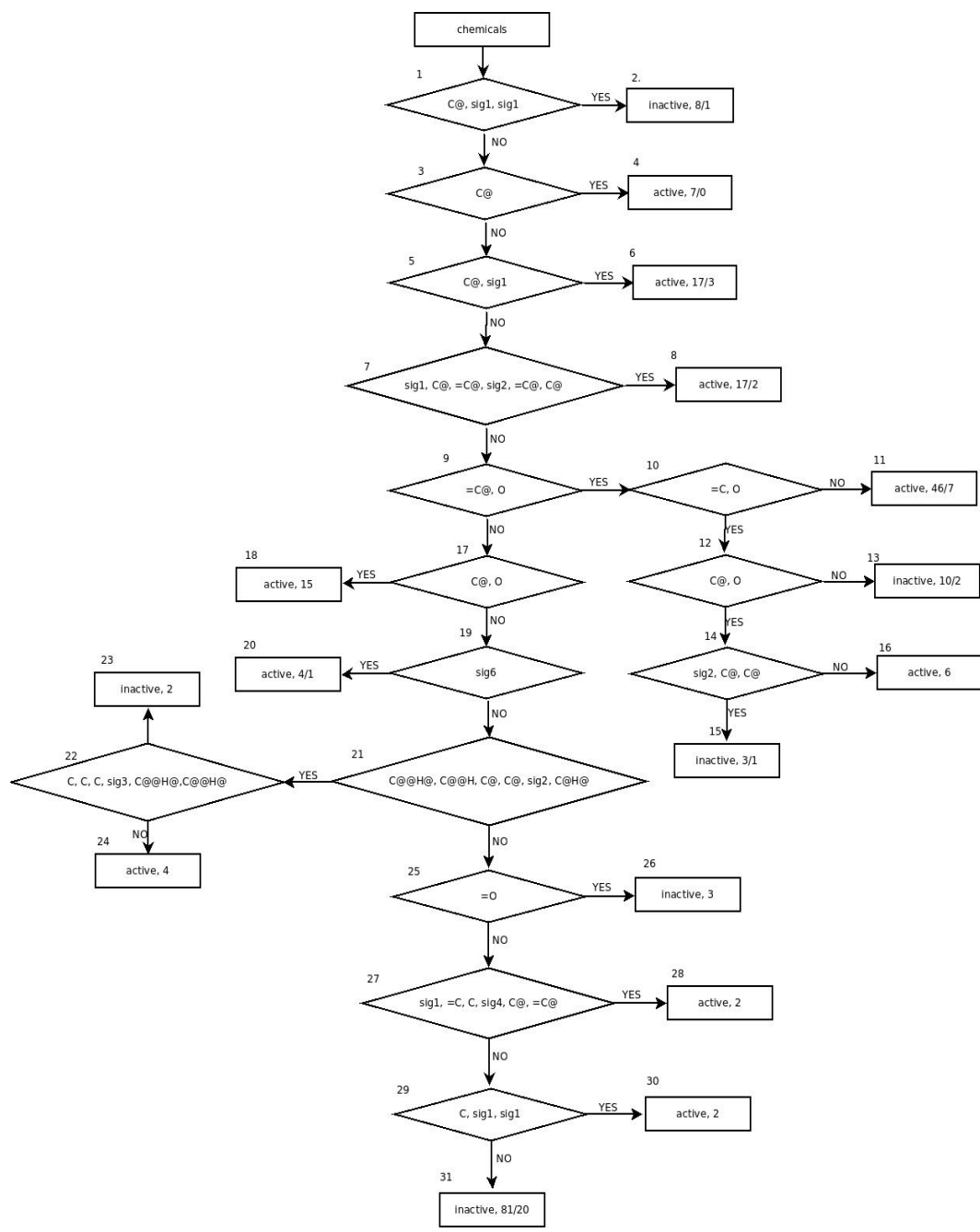
Figure 11



A molecule from the database:

O([Si](CC1C=CC=CC=1)(CC)[Si](CC2=CC=CC=C2)(C)C.

Figure 12



Method's predictive model in the form of IF- THEN- rules.