

1 **Towards a Partial Order Graph for Interactive Pharmacophore** 2 **Exploration: Extraction of Pharmacophores Activity Delta.**

3 Etienne Lehembre[‡], Johanna Giovannini^{*}, Damien Geslin^{*‡}, Alban Lepailleur^{*}, Jean-Luc
4 Lamotte[‡], David Auber^{*}, Abdelkader Ouali[‡], Bruno Cremilleux[‡], Albrecht Zimmermann[‡],
5 Bertrand Cuissart[‡] and Ronan Bureau^{*}

6 ^{*}Centre d'Etudes et de Recherche sur le Médicament de Normandie, Normandie Université,
7 UNICAEN, CERMN, 14000 Caen, France

8 [‡]Groupe de Recherche en Informatique, Image, Automatique et Instrumentation de Caen,
9 Normandie Université, UNICAEN, ENSICAEN, CNRS, GREYC, 14000 Caen, France

10 ^{*}Univ. Bordeaux, CNRS, Bordeaux INP, INRIA, LaBRI, Talence, France

11 **Abstract**

12 This paper presents a novel approach called Pharmacophore Activity Delta for extracting
13 outstanding pharmacophores from a chemogenomic dataset, with a specific focus on a kinase
14 target known as BCR-ABL. The method involves constructing a Hasse diagram, referred to as
15 the pharmacophore network, by utilizing the subgraph partial order as an initial step, leading to
16 the identification of pharmacophores for further evaluation. A pharmacophore is classified as a
17 'Pharmacophore Activity Delta' if its capability to effectively discriminate between active vs
18 inactive molecules significantly deviates (by at least δ standard deviations) from the mean
19 capability of its related pharmacophores. Among the 1,479 molecules associated to BCR-ABL
20 binding data, 130 Pharmacophore Activity Delta were identified. The pharmacophore network
21 reveals distinct regions associated with active and inactive molecules. The study includes a
22 discussion on representative key areas linked to different pharmacophores, emphasizing
23 structure-activity relationships.

24 **Keywords**

25 Hasse diagram, Partial order graph, Pharmacophore, BCR-ABL, Siblings, Activity Delta.

26

27 Introduction

28 The investigation of structure–activity relationships (Structure–Activity Relationships, SAR:
29 relationship between the structures of chemicals and their biological activities) represents one
30 of the most important tasks during the early stages of the drug discovery process [1]. The
31 definition of pharmacophores as a key to drug design is very well accepted in the field of
32 medicinal chemistry and is a key point to understand a molecule’s affinity for a biological
33 receptor [2]. In our initial publication on topological pharmacophores [3], we described the
34 logic for the definition of a new type of descriptor based on the notion of emergent
35 pharmacophores. We repeat some points here to clarify the objectives of this work.

36 A pharmacophore corresponds to the greatest common structural denominator associated with
37 a group of compounds exhibiting the same biological response profile [4]. Given a specific
38 target, ligand-based pharmacophore elucidation requires the detection of the spatial
39 arrangement of a combination of chemical features shared by several active molecules and
40 responsible for favorable interactions with the active site. To discover these common anchoring
41 features, the usual method starts with a careful selection of a small subset of ligands known for
42 binding to the same active site with the same binding mode [5]. We have done several studies
43 with this approach (see [6] for an example).

44 In recent years, the integration of large chemical databases [7] into the definition of SARs has
45 been clearly explored. With SARs and pharmacophores in mind, we have introduced a method
46 that automatically computes pharmacophores from a large data set of molecules without any
47 prior supervised selection of a small subset of molecules [3]. That method was based on the
48 computation of the so-called topological pharmacophores [8, 9].

49 Considering graph theory, 2D topological pharmacophores represent patterns which are present
50 in a number of chemical structures. When applied to a data set partitioned into two classes (*e.g.*,
51 active *vs.* inactive molecules), emerging pattern mining can identify the patterns that occur with
52 higher frequency in one of the two classes [3].

53 Of these topological pharmacophores, we can highlight those associated with particular
54 properties. We have previously explored a selection based on a growth rate value called GR
55 (Growth rate, GR: ratio of frequencies of appearance of a pharmacophore in a given class of
56 molecules compared to the other class (active or inactive compounds on BCR-ABL)). It
57 corresponds to the frequency of appearance of a pharmacophore in one class (active, for
58 example) compared to another. An initial selection was based on a value of 3 for the GR (ratio
59 of 3:1 for the frequencies between the two groups). A technique named Maximal Marginal
60 Relevance Feature Selection (Maximal Marginal Relevance Feature Selection, MMRFS :
61 selection of relevant pharmacophores by considering their number of associated chemicals and
62 their GR values) [10] has also allowed us to select a restricted subset of these topological
63 pharmacophores. This subset keeps the same statistical performance as the complete set
64 (sensitivity/specificity) with equivalent coverage of the compounds. First, pharmacophore
65 networks were defined based on these subsets by considering a graph editing distance [11] for
66 the calculation of the similarity between MMRFS pharmacophores and clustering techniques
67 [12]. For SAR studies and only for this objective, we have thought of inverting the frequencies
68 (inactive *vs.* active) and thus characterized topological pharmacophores associated with inactive
69 compounds. This gave us new insights into our data even if we are far from the historical
70 definition of pharmacophores.

71 In this study, we have chosen to focus on another view of our topological pharmacophores with
72 the definition of outstanding pharmacophores named Pharmacophore Activity Delta
73 (Pharmacophore Activity Delta, PAD : Pharmacophore for which the discrimination between
74 active *vs* inactive molecules significantly deviates from the mean capability of its related
75 pharmacophores.). To find these PADs, a Hasse diagram [13, 14] was defined as a
76 representation of the set of pharmacophores. This Hasse diagram corresponds to a partial order
77 graph [14, 15] encoding a partial order between pharmacophores, also called a pharmacophore
78 network. In this work, we leverage the pharmacophore network to quickly obtain the siblings
79 of a given pharmacophore.

80 For each pharmacophore, we quantify its level of significance using a quality measure function,
81 assigning a real number to each pharmacophore. We focus on the ratio of active molecules with
82 respect to a specific receptor, making the growth rate one of the functions used to assess the
83 quality of our pharmacophores. Pharmacophores that score very differently than the average of
84 neighboring pharmacophores are considered to be PADs. The definition of the neighbors is
85 based on the notion of siblings related to the Hasse diagram (*vide infra*). Using the growth rate,
86 we show experimentally that very few patterns turn out to be PADs.

87 **Methods**

88 **Dataset**

89 In line with our previous paper [12], we retrieved a ChEMBL compound data set of BCR-ABL
90 ligands [16–18] (target ChEMBL ID: ChEMBL1862, ChEMBL24 [19]). After discarding
91 compounds with molecular weight above or equal to 800 g/mol, we obtained a data set of 1479
92 molecules with either K_i or IC_{50} information. This limitation is primarily associated with the
93 combinatorial challenge when dealing with a molecule with a significant number of
94 pharmacophoric functions (*vide infra*). Of these 1479 molecules, 773 were designated active
95 compounds (meaning their K_i or IC_{50} value was below or equal to 100 nM).

96 **Pharmacophores**

97 In agreement with our previous description for the generation of pharmacophores [3, 12], the
98 pharmacophoric features correspond to generalized functionalities that are involved in
99 favorable interactions between ligands and targets, including hydrogen-bond acceptors ($|A|$)
100 and donors ($|D|$), negatively ($|N|$) and positively ($|P|$) charged ionizable groups, hydrophobic
101 regions ($|H|$), and aromatic rings ($|R|$). Therefore, a pharmacophore is a fully connected graph
102 where each vertex represents one of the specific pharmacophoric features, and the edges are
103 labeled with the number of the fewest possible bonds between two vertices. The number of
104 vertices, *i.e.* pharmacophoric features, composing a pharmacophore is called its order.

105 A notation was fixed for the pharmacophores. We started with the vertex of the
106 pharmacophores, *e.g.*, $|A|A|$ for a pharmacophore with two As with pipes as separators, and we
107 indicated the values of the edges, *e.g.*, $|2|$ for a distance of 2 bonds between the As, with pipes
108 as separators (final notation: $|A|A| |2|$). For a more complex case with four pharmacophoric
109 features and six distances to integrate, for instance $|A|A|H|D| |2|4|5|7|1|3|$, the six distances
110 correspond to the first one against the others ($|A|A|$, $|A|H|$, $|A|D|$) then, the second one against
111 the others ($|A|H|$, $|A|D|$) and, at the end, the last one against the other ($|H|D|$). In the following,
112 we omit edges information and “ $|$ ” separators in figures when they are not necessary for
113 comprehension.

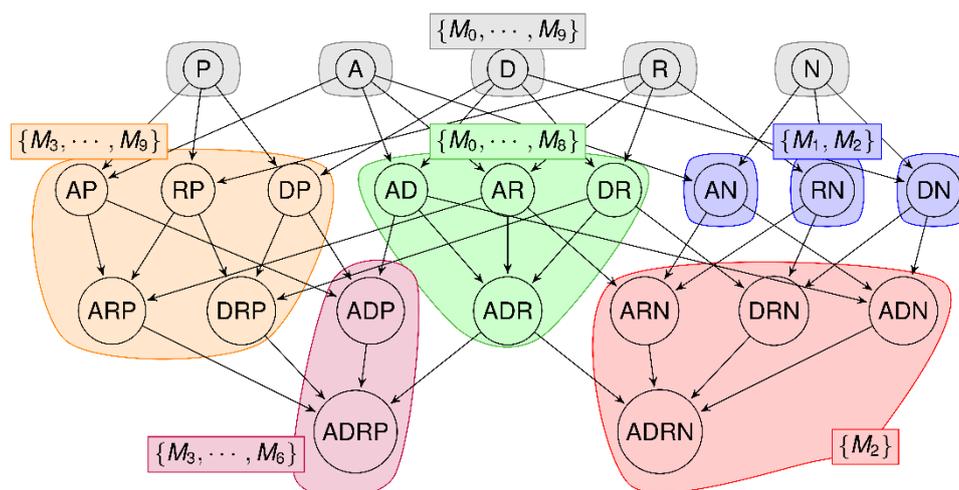
114 We call “support” the set of molecules supporting a given topological pharmacophore, *i.e.*,
115 containing all the pharmacophoric features of the pharmacophore with the correct distances
116 between them. Let p a pharmacophore and D the set of studied molecules. We note $Support(p)$
117 the support of p in D , *i.e.*, its set of supporting molecules.

118 In agreement with our previous studies [3, 12], the minimal support for the extraction of
119 pharmacophores was fixed to 10 (minimal number of compounds), and the orders (number of
120 pharmacophoric features) were between 1 and 7. 112 291 pharmacophores were generated with
121 these parameters. For the GR calculation (*vide infra*), the cutoff for active derivatives was fixed
122 to be less than or equal to 100 nM (773 compounds).

123 Let p, q be two 2D pharmacophores assimilated to graphs with labeled vertices and labeled
124 edges and D the set of studied molecules. If p is a subgraph of q , noted $p \subset q$, this means that
125 the pharmacophore p is included in the pharmacophore q . It also means that every single
126 molecule covered by q is also covered by p . A molecule set covered by a pharmacophore p is
127 the support of a pharmacophore denoted $Support(p) \subset D$. Thus, we can state that $p \subset q$
128 implies $Support(q) \subset Support(p)$.

129 From the subgraph partial order we can build a Hasse diagram[20] called a pharmacophore
130 network. We note $G(V, E)$ a pharmacophore network where each vertex $v \in V$ is a
131 pharmacophore and given two vertices $v_1, v_2 \in V, \exists (v_1, v_2) \in E$ (E is the set of edges between
132 the pharmacophore network vertices) if and only if $v_1 \subset v_2$ and $\nexists v_3 \in V$ such that $v_1 \subset v_3$
133 and $v_3 \subset v_2$. Therefore, an edge links two vertices of the network $v_1, v_2 \in V$ if and only if the
134 pharmacophore in v_1 is a subgraph of the pharmacophore contained in v_2 and there are no
135 pharmacophores v_3 in the pharmacophore network subgraph of v_2 which has v_1 as subgraph.

136 We note the edge relation between the vertices of the pharmacophore network $v_1 < v_2$ and call
137 v_1 a parent, which means that v_2 is called a child. It also means that $Support(v_2) \subset$
138 $Support(v_1)$. We illustrate the obtained structure in Figure 1.



139

140 **Figure 1.** Structure of the pharmacophore network. Each circle is a vertex containing a
141 pharmacophore. Only the pharmacophoric features are displayed to simplify the example and
142 the separators “|” are removed for ease of readability. Molecules having the pharmacophore are
143 indicated in the colored rectangles using set notation. The notation $\{M_3, \dots, M_6\}$ indicates that
144 the set is composed of molecules $M_3, M_4, M_5,$ and M_6 . The molecules associated to a
145 pharmacophore is determined by the colored area its vertex is in. Edges displays the inclusion
146 relation between pharmacophores. The vertex containing AN is connected to the vertices
147 containing ARN and ADN because AN is a subgraph of ARN and ADN. Since AN is associated
148 with molecules M_1 and M_2 , ARN and ADN must be associated to a subset of $\{M_1, M_2\}$. In this
149 example, these pharmacophores are associated to the molecule M_2 .

150 .

151 As we noticed that a large number of pharmacophores appear in the exact same set of molecules,
152 we decided to group them into equivalence classes [21] (ECs) based on molecule sets.

153 GEC, DEC, SEC

154 The first one is the General Equivalence Class (GEC), which groups every pharmacophore
155 covering the same set of molecules. Let p a pharmacophore and $G(V, E)$ a pharmacophore
156 network containing p , its general equivalence class is defined as $GEC(p, G) = \{v \in$
157 $V \mid Support(v) = Support(p)\}$. The formula can be transcribed as follows. Given a
158 pharmacophore p and a graph G , the general equivalence class of p is the set of pharmacophores
159 v contained in the vertices V of G having the same support as p , *i.e.* associated to the same set
160 of molecules. In **Figure 1**, these equivalence classes are indicated by the colors of the areas.
161 Meaning that pharmacophores of the first layer belong to the same general equivalence class
162 because they all are in grey areas.

163 The second one is the Divided Equivalence Class (DEC), which groups every pharmacophore
164 that has the same set of molecules and the same order. Let p a pharmacophore and $G(V, E)$ a
165 pharmacophore network containing p , we label $Order(p)$ its number of pharmacophoric
166 features. Then, the divided equivalence class of p is defined as $DEC(p, G) = \{v \in$
167 $GEC(p, G) \mid Order(v) = Order(p)\}$. The formula can be transcribed as follows. Given a
168 pharmacophore p and a graph G , the divided equivalence class of p is the set of
169 pharmacophores contained in the general equivalence class of p in G having the same
170 order, *i.e.*, having the same number of pharmacophoric features. In **Figure 1**,
171 pharmacophores in the orange area belong to the same general equivalence class but are
172 divided in two divided equivalence class regarding the layer they belong to, *i.e.*, regarding
173 their orders.

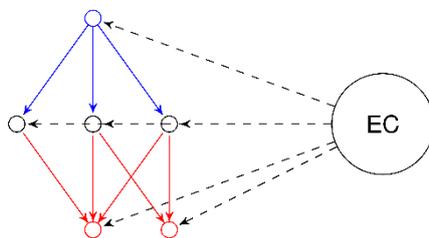
174 The last one is a specialization of GECs based on the connectivity of the pharmacophores in
175 the pharmacophore network called the Structured Equivalence Class (SEC). To define this
176 class, we introduce a new operator. Let p, v two pharmacophores in the vertices of the
177 pharmacophore network $G(V, E)$; we note $p \sim v$ if we have $p < v$ or $v < p$. Thus, given
178 $v_1, \dots, v_n \in V$, the expression $(p \sim v_1 \sim \dots \sim v_n \sim v)$ indicates that a path exists in the
179 pharmacophore network going from the vertex p to the vertex v . A structured equivalence class
180 groups all pharmacophores occurring in the same set of molecules having a path
181 connecting them inside their GEC. Let p a pharmacophore, its structured equivalence class
182 is defined as $SEC(p, G) = \{v \in GEC(p, G) \mid (p \sim v), \text{ or } (\exists v_1, \dots, v_n \in GEC(p, G), (p \sim$

183 $v_1 \sim \dots \sim v_n \sim v))\}$. The formula can be transcribed as follows. Given a pharmacophore
 184 p and a graph G , the structured equivalence class of p is the set of pharmacophores v in V
 185 contained in the general equivalence class of p in G which are connected to p by a path only
 186 visiting pharmacophores contained in the general equivalence class of p in G . In **Figure 1**, the
 187 pharmacophores in the grey areas all belong to the same general equivalence class but they all
 188 belong to separated structure equivalence classes.

189 The concepts of GEC, DEC and SEC all fall under a common concept called Equivalence
 190 Classes (EC). We can construct a pharmacophore network that minimizes redundant
 191 information from the ECs within a given pharmacophore network by taking ECs as vertices and
 192 extending the partial order as follows. Let EC_1, EC_2 be two equivalence classes; we say that
 193 $EC_1 < EC_2$ if and only if $\exists e_1 \in EC_1, \exists e_2 \in EC_2, e_1 < e_2$. With the extended partial order, we
 194 can define a pharmacophore network of equivalence classes following the same principles as
 195 the one used to compute the ECs. Below, we introduce methods applied to the pharmacophore
 196 network. These methods can be applied to either pharmacophores as vertices or equivalence
 197 classes as vertices. We refer to it as the GEC (respectively DEC and SEC) network when the
 198 vertices of the network are general (respectively divided and structured) equivalence classes.

199 In Figure 1, the three pharmacophores appearing only in the molecule M_1 and M_2 (in the blue
 200 areas) belongs to the same GEC and DEC, but do not belong to the same SEC because they are
 201 not connected within their GEC, *i.e.*, the path linking one to another in the context graph has to
 202 go through a vertex which covers different molecules set. But if you consider pharmacophores,
 203 ADN, DNR and ARN, they belong to the same SEC (the purple area) because ADRN is
 204 associated with the same set of molecules. As we can observe, the set inclusion of molecules is
 205 maintained, which indicates that there is an equivalence between the two types of
 206 pharmacophore network.

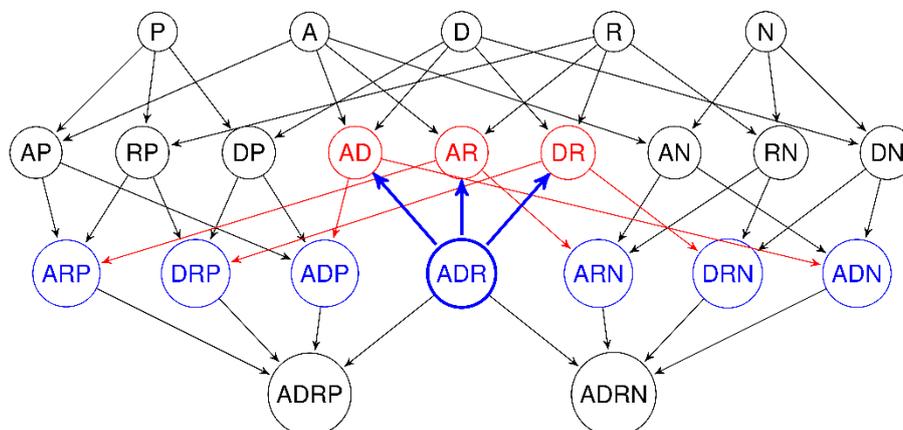
207 In order to study the equivalence classes, without considering every redundant pharmacophore
 208 contained, we use the notion of **generating pharmacophores** called **generators** and **closed**
 209 **pharmacophores**. Generators are pharmacophores that have no parents in their Equivalence
 210 Class (EC), which means they are the starting points of the EC. Closed pharmacophores are
 211 pharmacophores that have no children in their EC, which means they are the endpoints of their
 212 EC. In Figure 2, each circle without label is a pharmacophore (left) contained in the circle
 213 labeled EC which is an equivalence class (right). Dashed lines symbolize the inclusion relation
 214 between the pharmacophores and the equivalence class. We have one generator in blue and two
 215 closed pharmacophores in red.



216
 217 **Figure 2.** Generating pharmacophore (blue) and closed pharmacophores (red) from one EC
 218 (right).

219 But even with the use of equivalence classes, there are still too many vertices to study in the
 220 pharmacophore network. Therefore, we use the notion of siblings in a pharmacophore network.
 221 From intuition, a sibling is a vertex having at least one common parent. For a given

222 pharmacophore p and its pharmacophore network $G(V, E)$ where each vertex $v \in V$ is a
 223 pharmacophore, the siblings set of p is defined as $S(p, G) = \{v_1 \in V \mid \exists v_2 \in V, v_2 <$
 224 $p \text{ and } v_2 < v_1\}$. We note $Card(p, G)$ the cardinal of the set of siblings of p , *i.e.*, the number
 225 of pharmacophores contained in the siblings set.



226
 227 **Figure 3.** Getting the siblings ARP, DRP, ADP, ARN, DRN, and ADN (blue) from an origin
 228 vertex labeled ADR (bold blue); its parents are AD , AR , and DR (red).

229 In the pharmacophore network (see Figure 3), the siblings have at most one pharmacophoric
 230 feature which differs from the origin pharmacophore. In the condensed graph, the siblings cover
 231 the closest sets of molecules which are not included in one another because they all have
 232 molecule subsets of their common parents.

233 Using the concept of *siblings*, we will identify the *ECs* whose quality strongly deviates from
 234 those of their siblings. We interpret those *ECs* as key graph elements, as they may explain the
 235 biological behavior of their supporting molecules. We call in the following the selected
 236 outstanding pharmacophores the *Pharmacophore Activity Delta* (PAD).

237 Pharmacophore Activity Delta.

238 Let p a pharmacophore and D the molecule data set. We call the quality of p a real number
 239 determined by a function considering the molecules containing p noted $f(p, D)$. In this work,
 240 the quality is the normalized growth rate of p . We say that a pharmacophore p is a PAD when
 241 its quality deviates from the mean quality of its siblings $S(p, G)$. Let $f(p, D)$ the quality measure's
 242 value of the pharmacophore p in the dataset D , the sibling mean $\mu(S(p, G), D)$ is:

$$243 \quad \mu(S(p, G), D) = \frac{\sum_{s \in S(p, G)} f(s, D)}{Card(S(p, G))}$$

244 Equation 1

245 Then, $\sigma(S(p, G), D)$ is defined as the standard deviation of the siblings:

$$246 \quad \sigma(S(p, G), D) = \sqrt{\frac{\sum_{s \in S(p)} (f(s, D) - \mu(S(p, G), D))^2}{\text{Card}(S(p, G))}}$$

247 Equation 2

248 The pertinence of p is defined as:

$$249 \quad \text{Pert}(p, G, D) = \frac{f(p, D) - \mu(S(p, G), D)}{\sigma(S(p, G))}$$

250 The pertinence is the deviation from the mean quality of the sibling divided by the standard
251 deviation of the sibling. It can be transcribed as the deviation proportion of p regarding it
252 sibling. From this equation, a PAD is a pharmacophore which pertinence is high enough to
253 interest the expert. Therefore, we define our PAD selector.

254 The selector is defined as:

$$255 \quad \text{PAD}(G, f, D, \delta) = \{p \in V \mid |\text{Pert}(p, G, D)| \geq \delta\}$$

256 Equation 3

257 Thus, a pharmacophore p is a PAD if its quality deviates at least δ standard deviations (Equation
258 3) from the mean of the qualities of its siblings, δ being a user-supplied parameter.

259 We chose to use the standard deviation because we want to adapt our selection to each sibling.
260 If the siblings are different from one another, we only want to select the one that deviates the
261 most. If the siblings are similar to each other, then even a small deviation can be interesting.

262 GR

263 Based on the partitioning of the initial dataset into active and inactive molecules (or the inverse),
264 the growth rate (GR) of a given pharmacophore corresponds to the ratio between the frequencies
265 with which it occurs in each of the two subgroups.

$$266 \quad \text{GR} = \frac{\text{Fit frequency within actives}}{\text{Fit frequency within inactives}}$$

267 The main metric in this study is GR_N , normalized GR with values between 0 and 1.

$$268 \quad \text{GR}_N = \frac{\text{GR}}{\text{GR} + 1}$$

269 A GR value of 1 (same frequency for active and inactive compounds) corresponds to 0.5 for
270 GR_N . For the two extreme values, a GR_N value of 1 indicates that a pharmacophore occurs in
271 only active compounds, and a value of 0, in only inactive compounds. A GR value of 3,
272 classically used in our previous studies, now corresponds to a GR_N value of 0.75.

273 Pharmacophore (and PAD) stability

274 Discovering interesting substructures from data always risks capturing spurious phenomena
275 particular to the data set, instead of fundamental relationships that hold more generally. In the
276 case of pharmacophore activity deltas, this risk is compounded by the fact that each PADs
277 identification depends not only on its own support and quality, but also on those of its siblings
278 (and, furthermore, on whether those siblings are present in the pharmacophore network at all).

279 To assess the stability of discovered PADs, we therefore use a ten-fold cross-validation of the
280 data: the data set is split into ten equally-sized subsets (folds), which are then combined to
281 derive ten subsets, each of which containing 90% of the whole data, keeping one fold apart each
282 time. This allows to modify data sets in a controlled manner. PADs are identified independently
283 on each of those 10 data sets, and we assess how often PAD (re)occurs in the different result
284 sets.

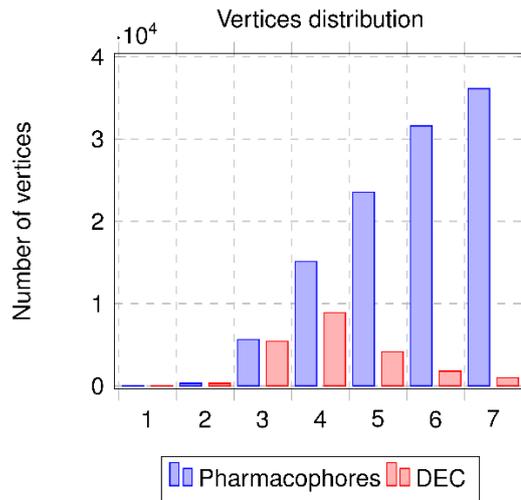
285 Given the construction of the underlying data sets, any two such sets will share a proportion of
286 about $1-10/90 = 0.88889$ of the compounds. Simply based on this data overlap, we would
287 expected particular pharmacophores to reoccur k times at most 0.88889^k due to chance (e.g.
288 0.5549 for $k=5$). As mentioned above, however, this probability will be significantly lower for
289 PADs since not only their siblings need to reoccur but GR differences will also need to be large
290 enough for a pharmacophore to be identified as a PAD.

291

292 Results and Discussion

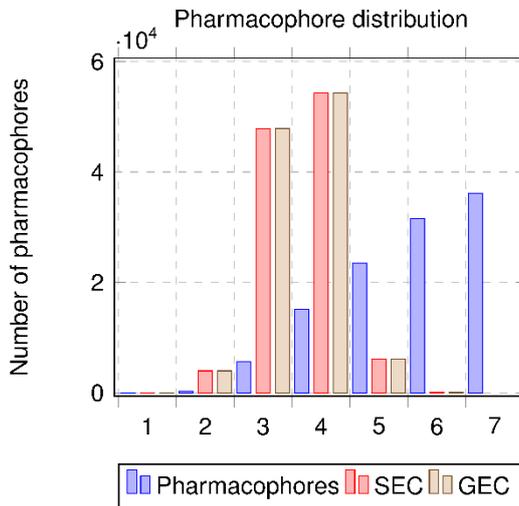
293 Pharmacophores and equivalence class network: GEC, DEC, SEC

294 Figure 4 shows the initial pharmacophore network (blue) and the DEC network (red) illustrating
295 the distribution of vertices regarding their layers. We can see that depending on the
296 pharmacophore order, the number of DEC vertices is strongly reduced when the order increases.
297 This phenomenon is predominant for orders 5, 6 and 7: those orders place a high number of
298 pharmacophores into an equivalence class when considering the DEC definition. A
299 multiplication of pharmacophores associated with the same set of molecules is clearly observed
300 and amplified when the number of pharmacophoric features is integrated.



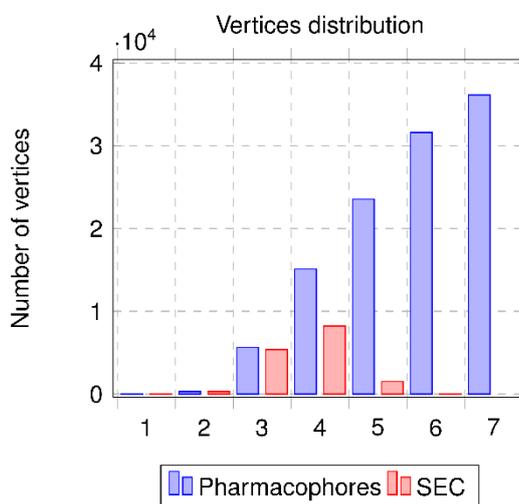
301
 302 **Figure 4.** Distributions of vertices by order: initial pharmacophore network and DEC network

303 Figure 5 shows the pharmacophore distributions for the initial pharmacophore network (blue),
 304 the GEC network (light brown) and the SEC network (red). For each EC, we have kept the
 305 number of initial pharmacophores for a particular view of the modifications. The new
 306 distribution of pharmacophores through the notion of ECs in the order is based on the generators
 307 (smallest pharmacophore for each EC) for each EC. We can see clearly that the GECs and SECs
 308 have the same distributions and the pharmacophores of orders 5–7 are redistributed, through
 309 the generators, to orders 3 and 4.



310
 311 **Figure 5.** Distributions of pharmacophores by order for initial pharmacophore network (blue)
 312 and for GECs network (light brown) and SECs network (red) when considering the generators
 313 for the distributions.

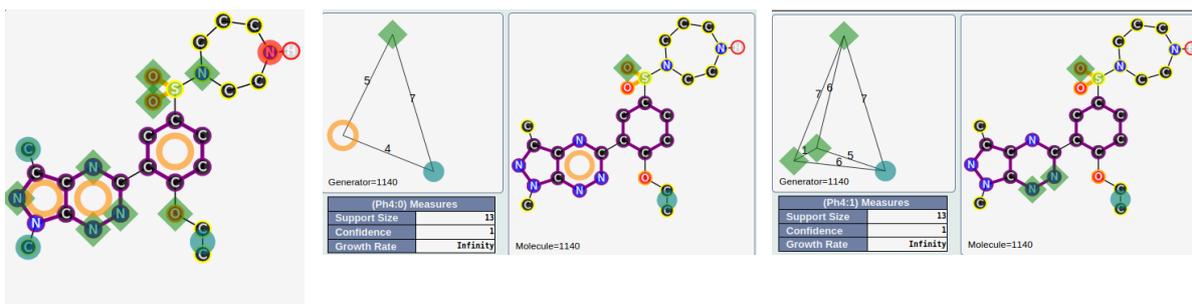
314 The last representation (see Figure 6) shows the distribution of vertices in the initial
315 pharmacophore network and in the SECs network by considering the order of the generators
316 for each SEC. From 112 291 pharmacophores in our initial data set, we move to 15477 SECs
317 to be assessed.



318
319 **Figure 6.** Distributions of vertices by order of pharmacophores for initial pharmacophore
320 network (blue) and SECs network (red) by considering the generators for the distributions and
321 one pharmacophore for each SEC.

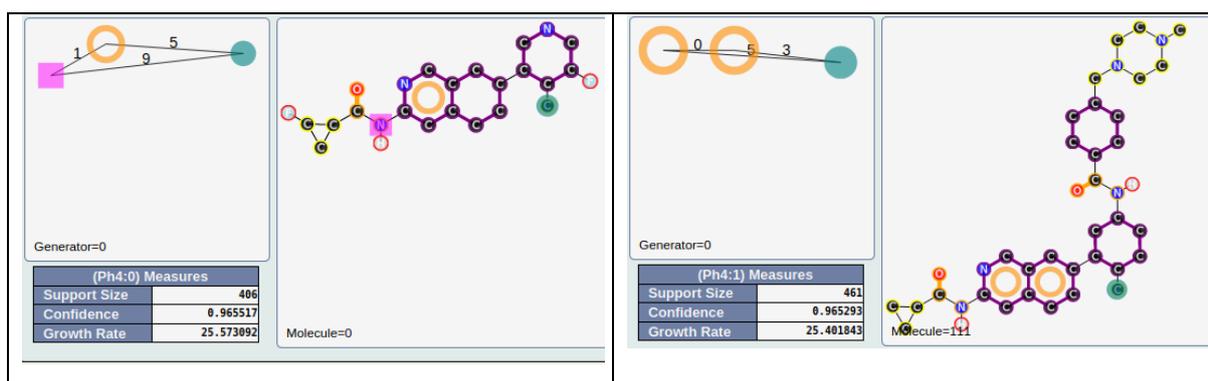
322 SEC/generators/parents

323 Of the 15477 SECs, 1745 are associated with at least two generators (*vide supra* for the
324 definition). These 1745 SECs cover 1301 out of 1479 compounds. Of these 1745 SECs, 443
325 are associated with at least 5 generators, 25 with at least 30 generators, 10 with at least 50
326 generators and 1 with 271 generators. To give a first explanation of these results, the size of the
327 molecules and the associated number of pharmacophoric functions were analyzed. In the initial
328 dataset, 99 compounds have a molecular weight ≥ 500 g/mol and a number of pharmacophoric
329 functions ≥ 20 . Of these 99 compounds, 51 are associated with a SEC with at least 30
330 generators. So, the molecular weight and the number of associated pharmacophoric functions
331 give a first and clear explanation for the observed number of generators associated with some
332 SECs. The SEC with 271 generators corresponds to 13 compounds, all inactive (see Figure 7
333 for illustrations of this SEC), compounds in agreement with the previous remark.



334 **Figure 7.** All pharmacophoric functions associated with one representative compound (left).
 335 Two among the 271 generators of this SEC (center and right).

336 Starting from the 1745 previous SECs, we analyzed the parents of these SECs. We wanted to
 337 see if some parents have particular characteristics in terms of filiations. These 1745 SECs are
 338 associated with 7517 parents. Among them, 669 have more than 20 filiations and 201 have
 339 more than 50 filiations. Among the last group, 18 parents are associated with active compounds
 340 ($GR_N \geq 0.75$) and 5 parents have a GR_N value ≥ 0.9 . The best ones (w.r.t. GR_N values), |R|D|H|
 341 |1|5|9| and |R|R|H| |0|3|5|, are associated with 406 and 461 compounds, respectively (see Figure
 342 8). These pharmacophores correspond to important pharmacophores of this kinase with
 343 structural characteristics associated with the interaction with the hinge region and the key
 344 methyl group, often related to the back pocket region of the binding site (see Xing et al. [22]
 345 and our latest publication [12]).



346 **Figure 8.** Best parents with more than 50 filiations.

347 Among the parents, the best one for the number of filiations, |A|R| |2|, has 276 filiations. It
 348 covers 1366 compounds out of 1475, with a GR_N value of 0.53.

349 Outstanding pharmacophores: from SECs to PADs

350 With EP mining in mind, we applied the SEC network to our pharmacophore file and retrieved
 351 the GR_N values for each SEC. Table 1 shows information about the distribution of SECs as a
 352 function of the GR_N values. For the selection of PADs among the SECs, the pertinence value
 353 (Equation 3, δ value) was considered first to be 1.96 (p-value of 0.05). 42 PADs (Table 1) were
 354 obtained, but with low coverage of the initial data set (22%). As a result, we have lowered the
 355 pertinence value to 1.64 (p-value of 0.1), and 377 PADs were obtained with a coverage of 81%
 356 of the initial data set (of the 277 compounds missing, 75 are actives).

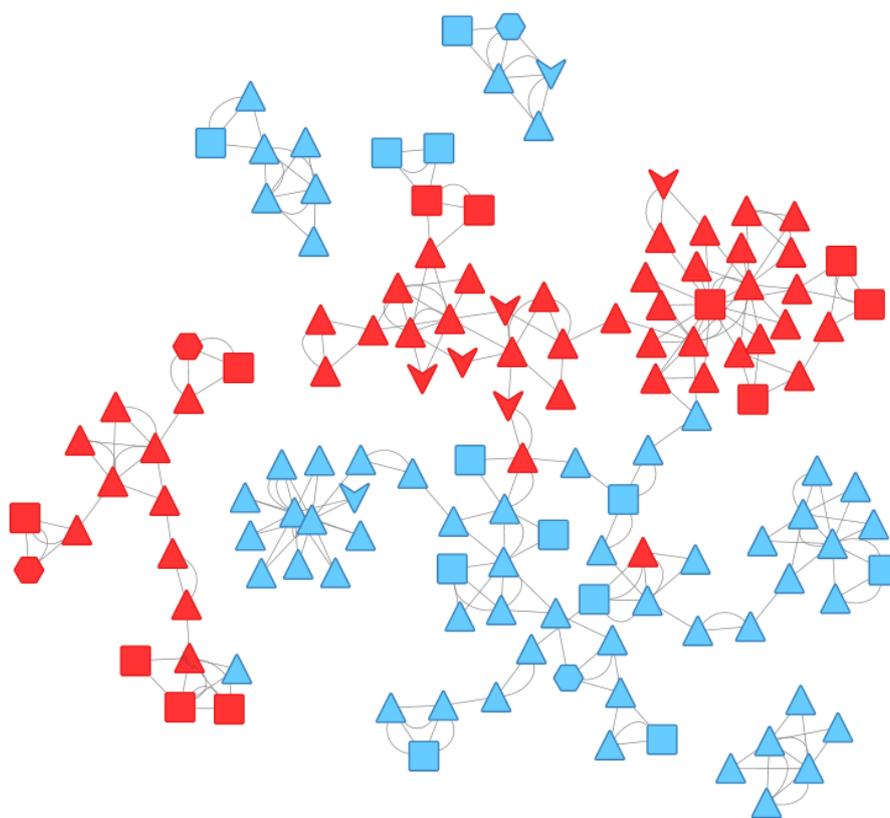
357 To analyze the PADs, we have chosen to represent them as a pharmacophore network. A
 358 similarity matrix was defined for the initial chemical data set with ECFP4 as molecular
 359 fingerprint descriptors. The Tanimoto coefficient was used as the similarity measure. The

360 similarity between the PADs was defined as the average similarity between the molecules
361 associated with the PADs. The orders of the PADs are different, so it was impossible to integrate
362 a graph edit distance in agreement with our previous studies for the similarities between the
363 PADs [12]. To decrease the number of PADs and in line with our initial studies, we decided to
364 summarize the initial PADs set using the MMRFS technic. The method is described in a
365 previous publication [3]. MMRFS aims to generate a subset of pharmacophores characterized
366 by discriminating, distinct, and representative elements of the active molecules. 135 PADs (see
367 Table 1) out of the 377 initial PADs were selected in this case with a coverage, for the data set,
368 of 77% (instead of 81% without MMRFS selection). Most of the PADs have order 3,
369 corresponding to three pharmacophoric functions (see Table 1). As described in our previous
370 publication, we focused, for the network, on the nearest neighbors of each PAD. The neighbors
371 of each PAD were ranked in descending order based on similarity coefficient values. Using this
372 method, the nearest two neighbors of each pharmacophore were retained (we also analyzed the
373 nearest five and ten neighbors, but the nearest two neighbors were best for the analysis of the
374 network). The two neighbors corresponded to the minimum number of neighbors because
375 several of the edges within a given network can exhibit identical values for the similarity
376 coefficient. We have chosen the Compound Spring Embedder[23] for the layout (PAD network)
377 in Cytoscape [24]. The final PAD network allows us to get a view of our data set (see Figure 9)
378 with active PADs in solid red and inactive PADs in solid cyan.

379 **Table 1.** Description of SECs (as a function of GR_N values) and PADs (as a function of
 380 pertinence values). Information on the associated number of SECs or PADs and the molecules
 381 covered for the PADs. Pertinence is related the Equation 3.

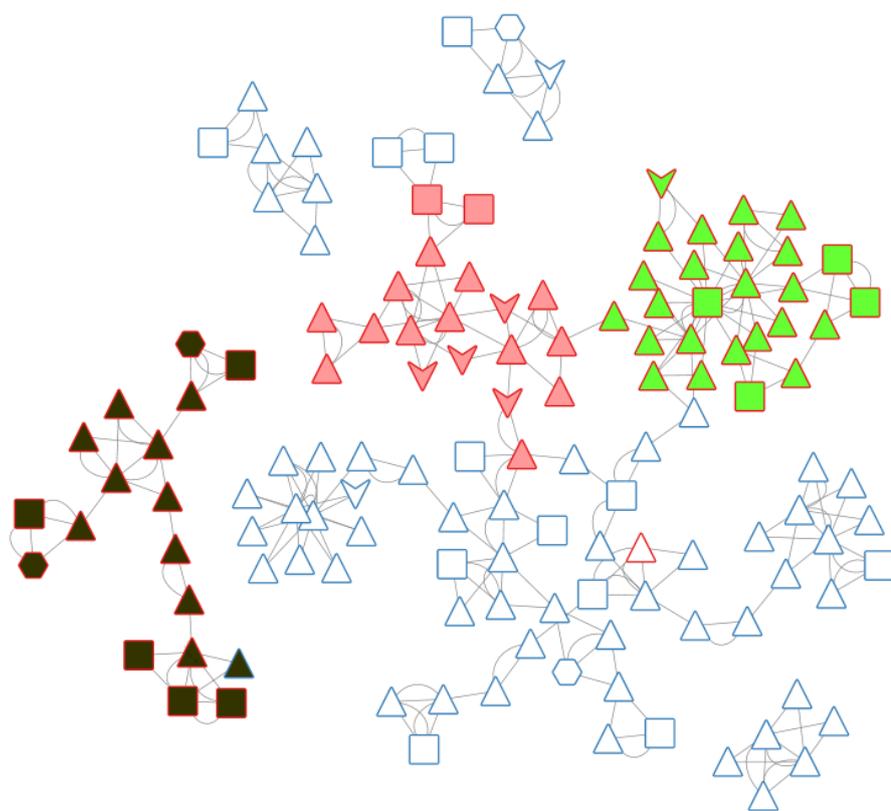
Descriptions of SECs	Number
SECs: $GR_N \geq 0.5 / GR_N < 0.5$	7534/7926 (SECs)
SECs: $GR_N \geq 0.75 / GR_N \leq 0.25$	4285/3803 (SECs)
SECs: $GR_N = 1 / GR_N = 0$	1084/979 (SECs)
Description of PADs	Number
PADs: pertinence ≥ 1.96 or ≤ -1.96 ($\alpha = 0.05$)	20/22 (PADs)
Molecules covered (active/inactive)	337 (molecules) (187/150)
PADs: pertinence $\geq 1,64$ or ≤ -1.64 ($\alpha = 0.1$)	187/190 (PADs)
Molecules covered (active/inactive)	1202 (molecules) (698/504)
PADs MMRFS with $\alpha = 0.1$	63 (0.85)/72 (0.60) (PADs (Recall values))
Molecules covered (active/inactive for above PADs MMRFS (// as separator))	659/44 // 19/426 (molecules)
Order 2 PADs	5/2 (PADs)
Order 3 PADs	45/56 (PADs)
Order 4 PADs	11/12 (PADs)
Order 5 PADs	2/2 (PADs)

382



383
 384 **Figure 9.** Pharmacophore Activity Delta network with active pharmacophores in red and
 385 inactive pharmacophores in cyan. The symbols are related to the order of the pharmacophores
 386 (order 2 (arrow), order 3 (triangle), order 4 (square), order 5 (hexagon))).

387 From this PAD network, we can distinguish, at first glance, three areas for active PADs (see
 388 Figure 10). A description of the PADs (in brackets, the number of associated chemicals) for
 389 these three areas and their representative compounds are provided in Tables 2, 3 and 4. The first
 390 one, in solid green, groups 26 PADs and covers 467 molecules (442 actives) with a recall value
 391 of 0.57 (57% of all actives). The second one, in solid pink, groups 19 PADs and covers 179
 392 molecules (170 actives). The last area, in solid black, is more isolated. It groups 18 PADs, 17
 393 actives and one inactive. The 17 active PADs cover 175 molecules (160 actives).

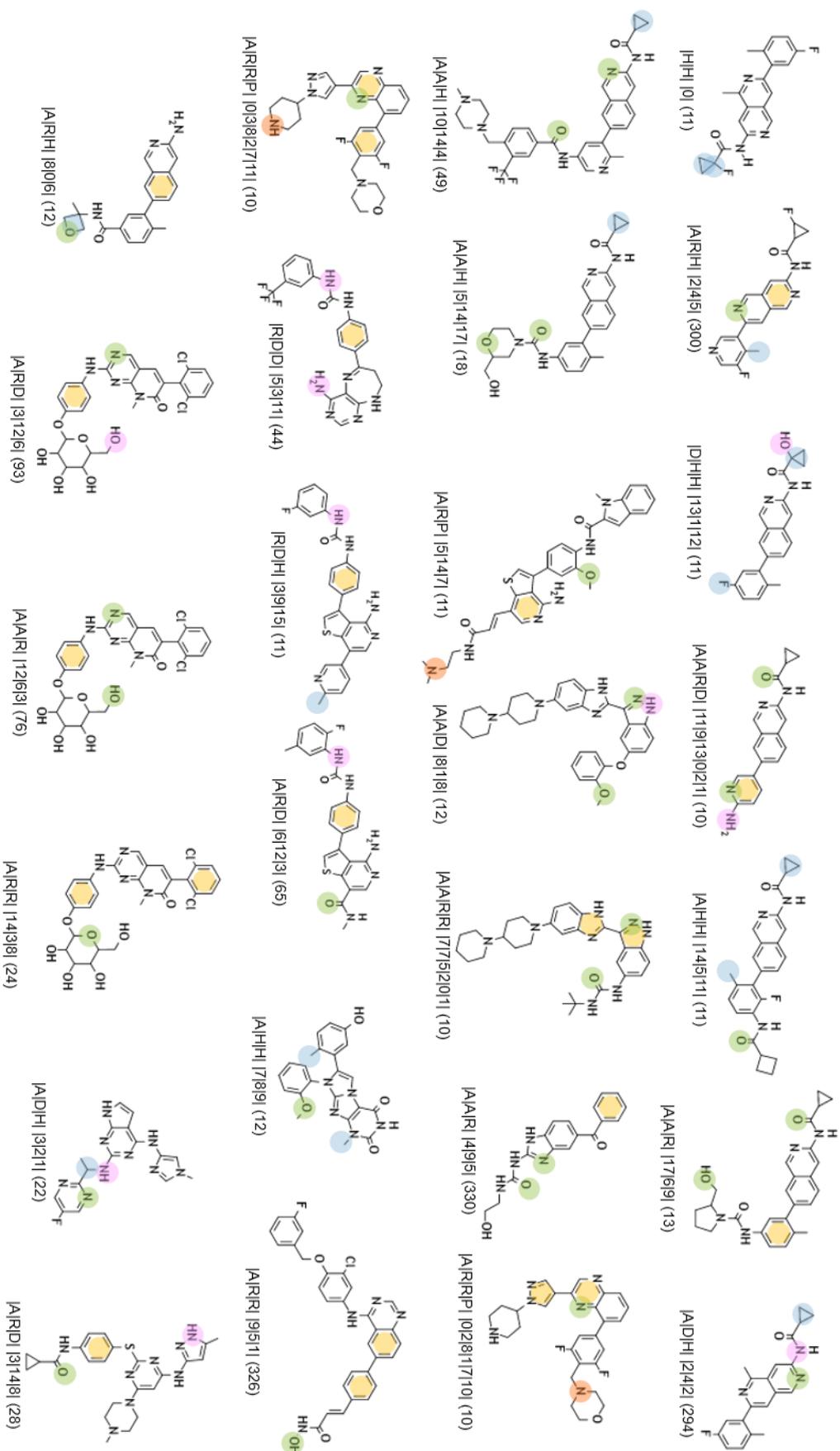


394

395 **Figure 10.** PAD network with the three areas (green, pink and black) for active
 396 pharmacophores.

397

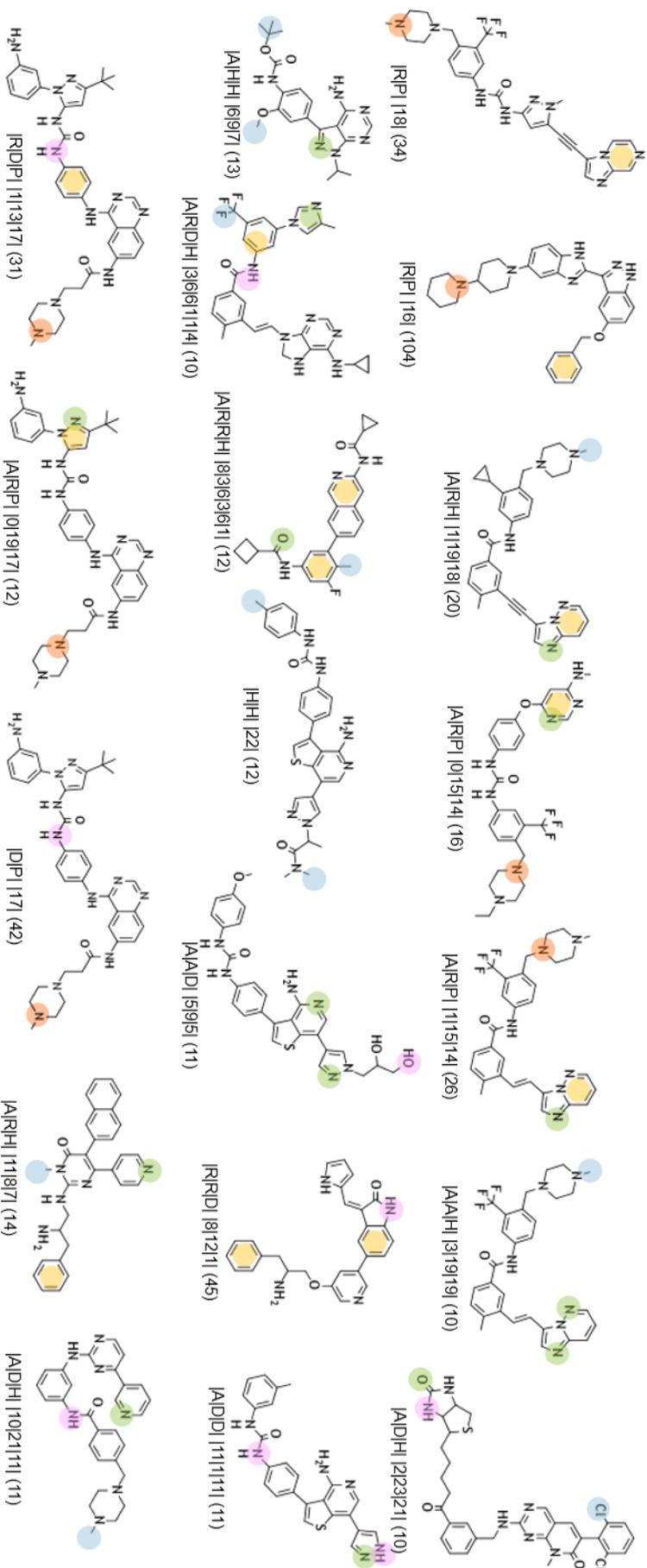
398 **Table 2.** Description of the representative compounds (centroid, ECFP4/Tanimoto) associated
 399 to each pharmacophore (with the alignment between molecules and pharmacophores) in the
 400 green area and in brackets, the number of associated chemicals.



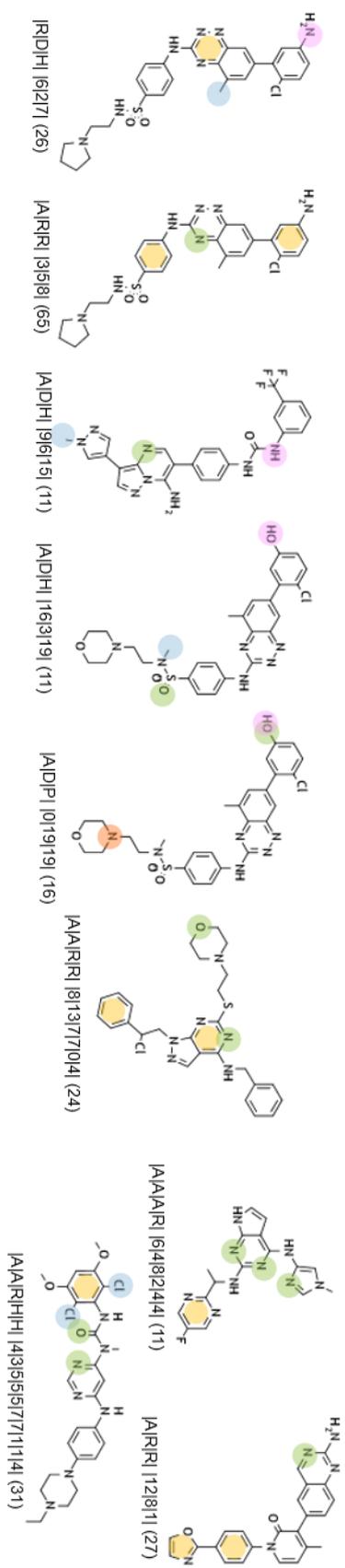
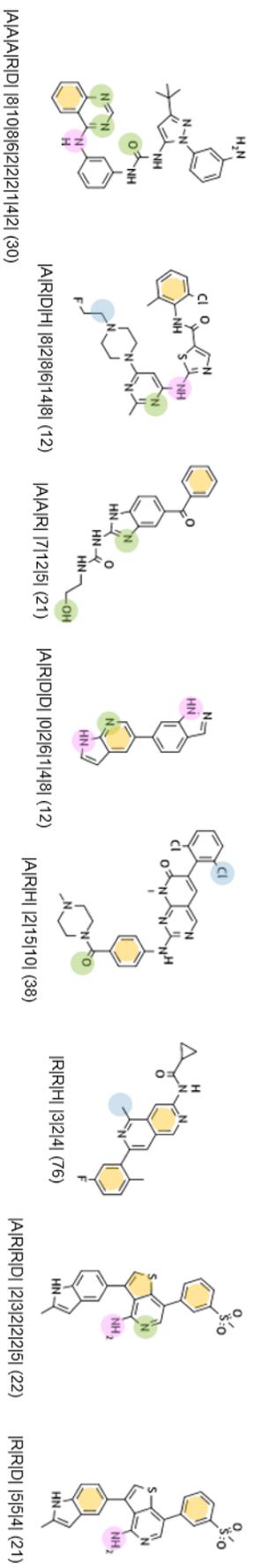
Pharmacophore color-code: ● Hydrogen-bond acceptor [A] ● Hydrogen-bond donor [D] ● Hydrophobic region [H] ● Positively charged/ionizable group [P] ● Aromatic ring [R]

402 **Table 3.** Description of the representative compounds (centroid, ECFP4/Tanimoto) associated
403 to each pharmacophore (with the alignment between molecules and pharmacophores) in the
404 pink area and in brackets, the number of associated chemicals.

Pharmacophore color-code: ● Hydrogen-bond acceptor [A] ● Hydrogen-bond donor [D] ● Hydrophobic region [H] ● Positively charged ionizable group [P] ● Aromatic ring [R]



406 **Table 4.** Description of the representative compounds (centroid, ECFP4/Tanimoto) associated
407 to each pharmacophore (with the alignment between molecules and pharmacophores) in the
408 black area and in brackets, the number of associated chemicals.



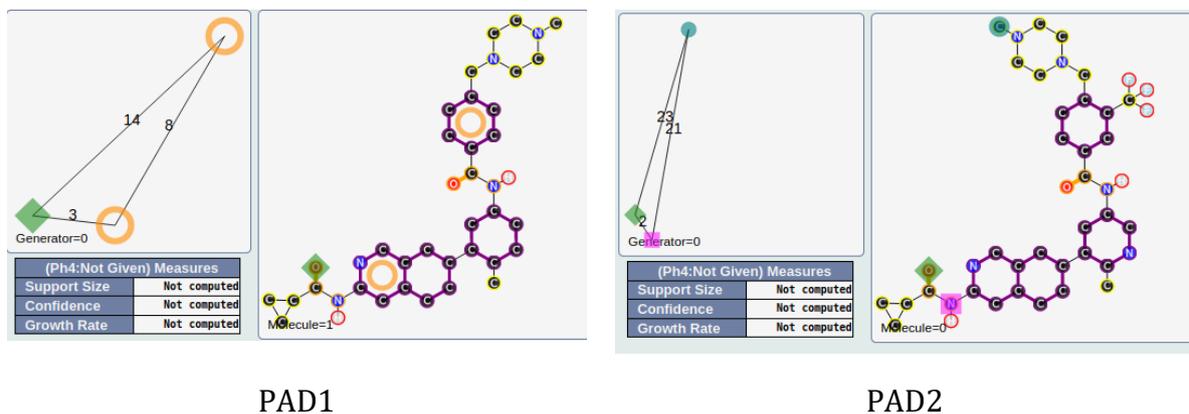
Pharmacophore color-code: ● Hydrogen-bond acceptor [A] ● Hydrogen-bond donor [D] ● Hydrophobic region [H] ● Positively charged ionizable group [P] ● Aromatic ring [R]

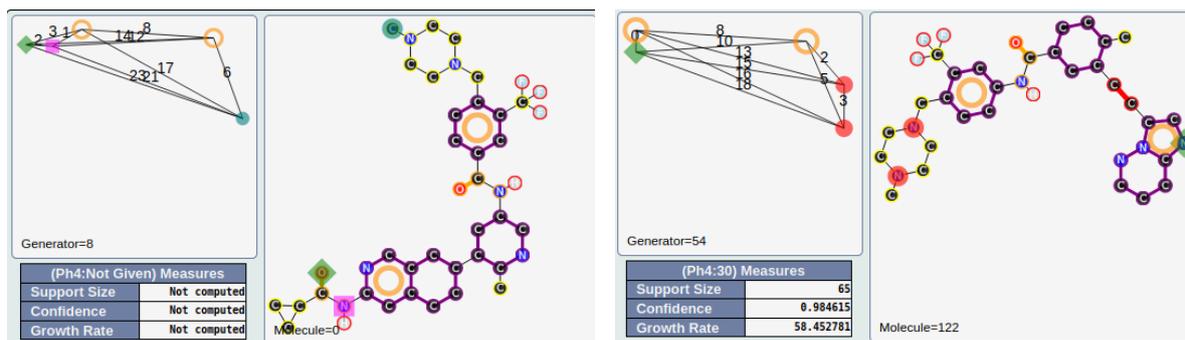
410 It is impossible in this publication to describe all the PADs. We have therefore chosen to focus
 411 on some specific areas of this PAD network. The first one concerns a connection between two
 412 active areas. In fact, the green area is connected to the pink area by two PADs (similarity of
 413 0.39 between the two PADs; see Figure 10 and Figure 11 (left), solid red).



414 **Figure 11.** PADs between two active groups (left, solid red) and proximities between active
 415 and inactive pharmacophores (right) corresponding to three areas (solid yellow, green, blue).

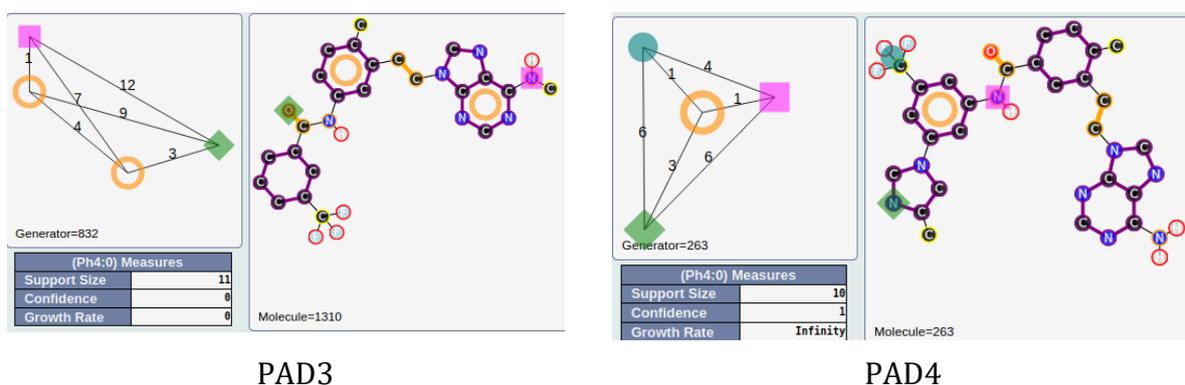
416 One of these two PADs has 24 compounds (PAD1, |A|R|R| |14|3|8|; see Figure 12), and the other
 417 has 10 compounds (PAD2 |A|D|H| |2|23|21|; see Figure 12). They share 9 compounds (90% of
 418 the compounds associated with PAD2 are in PAD1). By combining PAD1 and PAD2,
 419 pharmacophore **1**, with 5 pharmacophoric features, can be derived (see Figure 12). To analyze
 420 the possible proximity of other scaffolds to these nine compounds, we derived all the
 421 pharmacophores with 5 features from these nine compounds and extracted those associated with
 422 the maximum number of derivatives. Among the 56 pharmacophores generated, the best one
 423 (in terms of the number of compounds) is associated with 65 chemicals (pharmacophore **2**, GR
 424 = 58) and is related to the ponatinib-like family [25].





425 **Figure 12.** PAD1 and PAD2 with two representative compounds (top). Below, pharmacophore
 426 **1** corresponding to a combination of PAD1 and PAD2 (left) and pharmacophore **2** (with
 427 ponatinib) derived from the nine compounds fitting pharmacophore **1**.

428 For the other areas, we analyzed the situation where active PADs are close to inactive ones
 429 (similarity ≥ 0.3 for the PADs). This is the case for three areas (solid yellow, green, blue; see
 430 Figure 11). The first one, in solid yellow, allows us to understand the importance of one |A|
 431 function included in an aromatic group and a specific position of the |D| function for similar
 432 compounds. In fact, for PAD3 with only inactive compounds, we observed (see Figure 13) for
 433 the representative compound the inversion of the amide function compared to the representative
 434 compounds of PAD4 (with only active compounds), and moreover, one |A| function is missing
 435 compared to the representative compound of PAD4. PAD4 is the pharmacophore clearly
 436 associated with the nilotinib-like family [26].

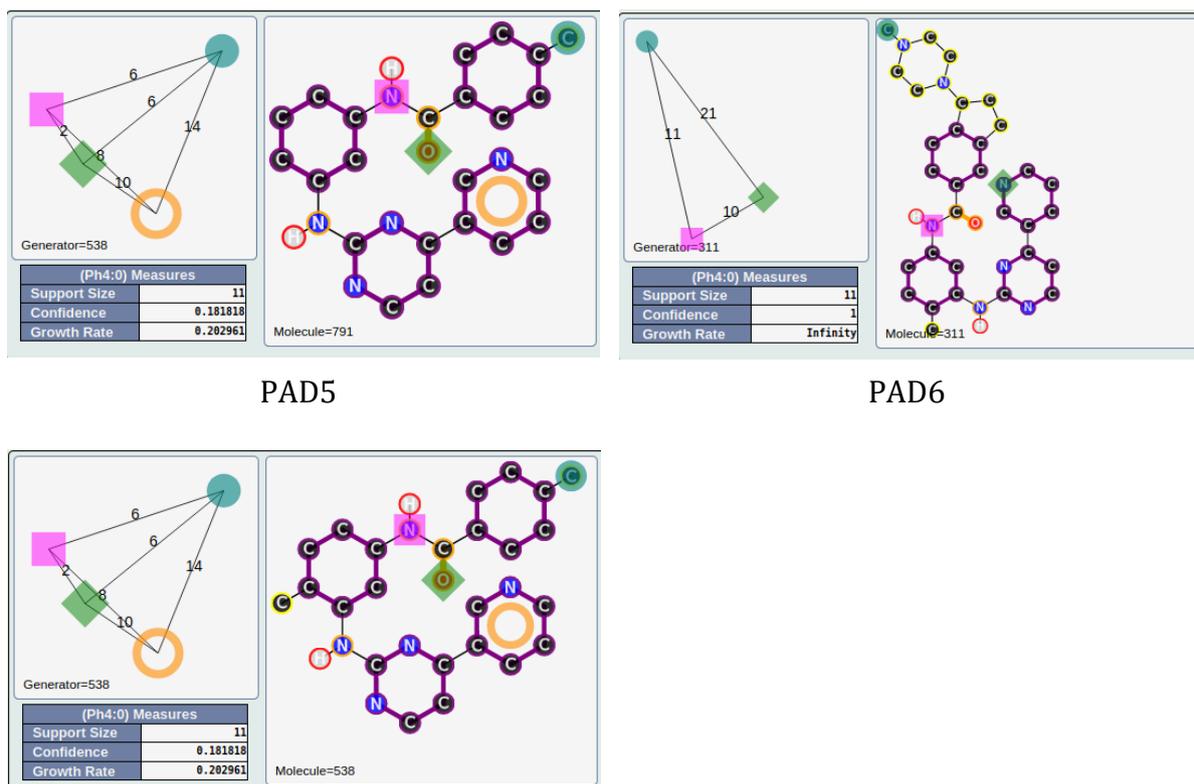


PAD3

PAD4

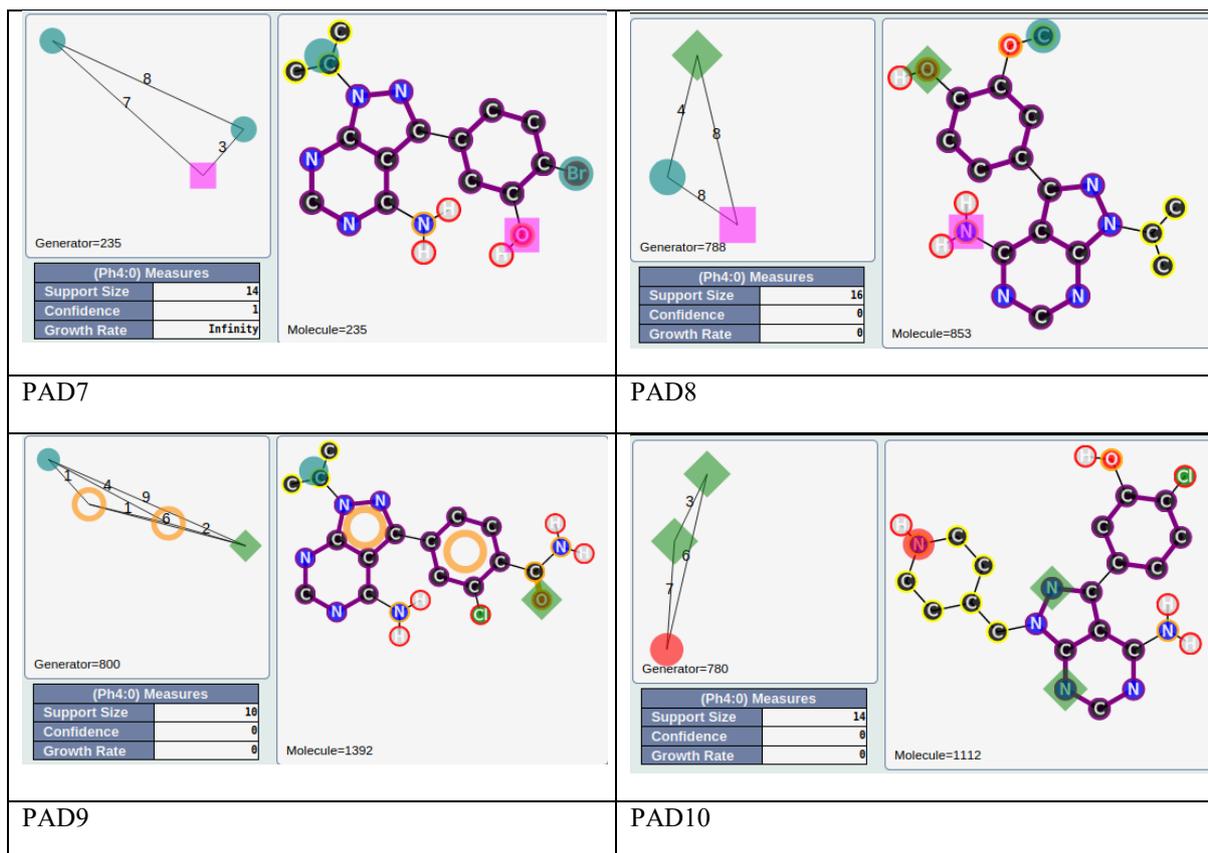
437 **Figure 13.** PAD3 (only inactive compounds) and PAD4 (only active compounds) showing the
 438 inversion of the amide function (between two aromatic rings) for the two representative
 439 compounds of each PAD.

440 The solid green area is associated with pharmacophoric variations around the scaffold
 441 associated with imatinib [27]. PAD6 (see Figure 14) has three pharmacophoric functions
 442 translating the position of two polar functions (|A| and |D|) and, above all, the size of the
 443 compound with a hydrophobic group being at a distance of 19 (19 edges) from the aromatic
 444 ring bearing the |A| function. PAD5 is, on the contrary, associated with inactive derivatives. We
 445 observed with PAD5 the typical scaffold associated with imatinib without the terminal amine
 446 functions. We can notice that one compound fitting PAD5 is active with the typical methyl
 447 group for some kinase inhibitors of ABL1 related to the back pocket binding site,[12] also
 448 described previously with the best parent (see Figure 8).



449 **Figure 14.** PAD5 and PAD6 with two representative compounds (top). PAD5 with the only
 450 active compound (bottom) associated to this PAD and for which the key methyl group is
 451 present.

452 The last solid blue area is associated with the only active PAD surrounded by inactive PADs
 453 (see Figure 9). PAD7 is active, and the other ones are inactive (see Figure 15). For this chemical
 454 series, we can clearly see the importance of the |H| function in alpha position to the |A|
 455 of the phenol group (PAD7). For PAD8, always on the phenol group, the |H| function is not in
 456 the same position (methoxy group in this case). For PAD 9, we do not have a phenol group and
 457 the |A| function is in a different position. For PAD10, a |P| function is present. So, some clear
 458 structure–activity relationships could be identified from the analysis of these PADs.



459 **Figure 15.** PAD7, PAD8, PAD9, PAD10 with representative compounds. PAD7 is active and
 460 the other ones are inactive. PAD7 and PAD8 have different positions of the hydroxy group for
 461 phenol functions. No phenol group in PAD9 compared to PAD7. PAD10 with a polar function
 462 (amine) instead of a hydrophobic function for PAD7.

463 Cross-validation studies and stability of PADs

464 We performed a stratified 10-fold cross-validation study (i.e. each fold contained the same
 465 proportions of active and inactive compounds) on the initial dataset (using scikit-learn/Kfold
 466 [28]). The method (extraction of pharmacophores and definitions of pharmacophore
 467 network/SECs/PADs) was applied independently to each subset.

468 As implied by, and supporting, our earlier explanation, we extract on average 15.6% fewer
 469 pharmacophores (5.7%-31.1%). The number of SECs varies less, between 7.75% and 13.5%
 470 fewer, for an average of 10.9% fewer SECs. A total of 364 active PADs were obtained by
 471 combining the result derived from the 10 subsets, fewer than the 377 PADs derived from the
 472 full data. The method was found to be more stable than we expected. Indeed, of these active
 473 PADs, 61% are present in at least 5 subsets (as compared to the 55.49% of *pharmacophores*
 474 one would expect to reoccur 5 times) and 26 PADs are present in the results of *all* the subsets.
 475 Of these 26 active PADs, |D|A|R| |2|3|2|, with 470 compounds, is associated with the highest
 476 number of compounds. Inactive PADs are less stable, with 40% present in at least 5 folds, and
 477 14 PADs are present in all the folds.

478 **Table 5.** Number of pharmacophores (Phar.), SEC and PADs for each fold. Cumulative
 479 presence of the PADs in the folds.

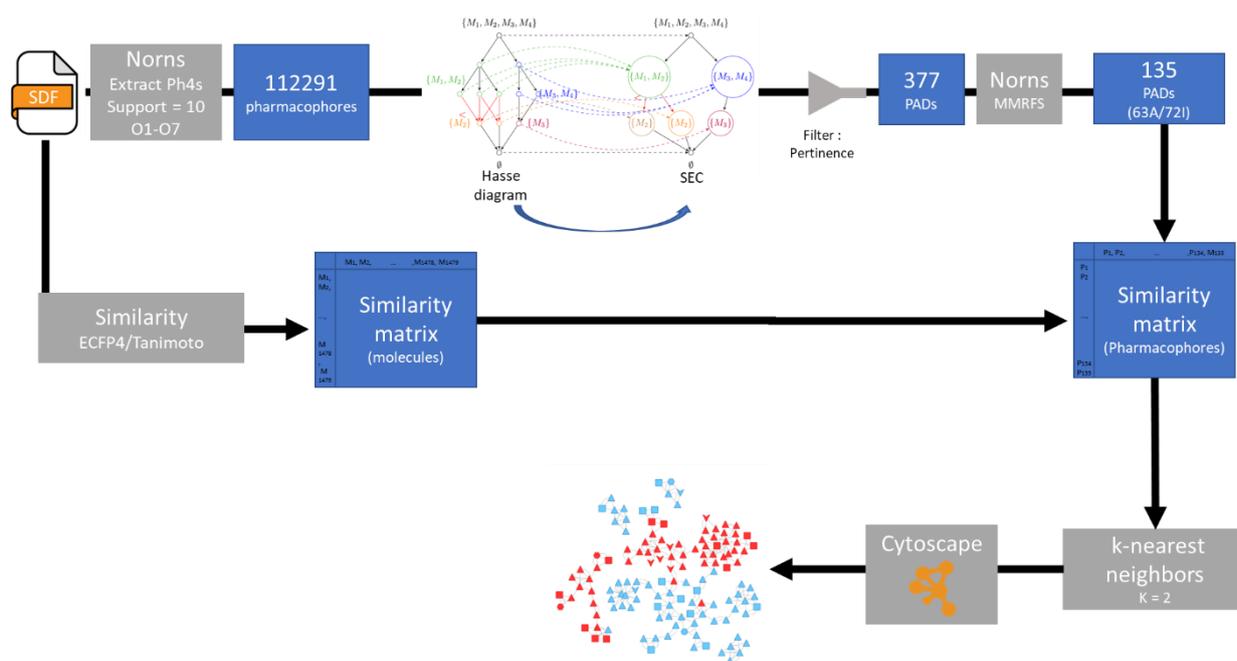
	F1	F2	F3	F4	F5	F6	F7	F8	F9	F10
Phar.	98053	92727	77346	92742	105875	95596	108556	90708	91588	93638
SEC	14257	14278	13836	13717	14058	13529	13378	13415	13547	13791
PADs	153/ 199	168/ 173	174/ 144	152/ 122	130/ 126	117/ 149	126/ 148	130/ 208	127/ 180	115/ 169
Presence of the PADs in the folds (cumulative: at least x folds)										
x fold	10	9	8	7	6	5	4	3	2	1
Pertinence ≥ 1.64	26	70	111	146	192	225	261	304	349	364
Pertinence ≤ -1.64	14	33	60	101	139	194	221	268	355	481

480 Conclusions

481 In an effort to develop a tool that can rapidly provide information from a dataset of molecules
482 regarding active or inactive compound characteristics, we conducted structural elucidation
483 using a fully annotated dataset of molecules extracted from the ChEMBL database. The various
484 steps involved in this workflow are summarized in **Figure 16**. The extraction of
485 pharmacophores with Norns is the most time-consuming process, taking several minutes with
486 our configuration. We processed 1479 molecules to generate topological pharmacophores
487 containing 1 to 7 motifs, with the support of at least 10 molecules. As part of our objective to
488 involve a human expert in pharmacophore elucidation, we established a specific method to
489 identify outstanding pharmacophores known as PADs.

490 The extraction of PADs is initially linked to defining the 15477 SECs from the initial 112291
491 pharmacophores. Subsequently, calculations of GR_N were performed for each SEC. A threshold
492 for the pertinence values associated with each SEC led to the extraction of 377 PADs. In the
493 end, a PAD network was constructed using Cytoscape starting from a representative set of 135
494 PADs (MMRFS). This network incorporates the similarity between the PADs for link definition
495 (the 2NNs of each PAD).

496 The interestingness of this reduced set of 135 PADs is based on the diversity of information it
497 provided, equally shared between active and inactive compounds. Cross-validation studies can
498 be also a basis for the selection of interesting PADs. The proximity between these PADs allows
499 us to explain some key SARs with the four illustrations in this publication.



500

501 **Figure 16.** Workflow associated to the different steps of our process from the extraction of
 502 pharmacophores to the representation of a network associated to the PADS.

503

504 **List of abbreviations**

505 **EP:** Emerging Pattern

506 **GR:** Growth Rate

507 **MMRFS:** Maximal Marginal Relevance Feature Selection

508 **EC:** Equivalent Class

509 **GEC:** General Equivalent Class

510 **DEC:** Divided Equivalent Class

511 **SEC:** Structured Equivalent Class

512 **PAD:** Pharmacophore Activity Delta

513 **SAR:** Structure-Activity Relationships.

514 **Declarations**

515 **Availability of data and materials**

516 The datasets supporting the conclusions of this article and the main programs related to this
 517 work are available at https://osf.io/pj8n3/?view_only=f49d5a0af5114b568f327c11d46bdfd3.

518 The main program and the sources (<https://hal.science/hal-04057516>) are available on GitHub:
 519 <https://github.com/Etienne-Lehembre/Pharmacophores-Activity-Delta>.

520 Image_Dockers for Norns are available on docker hub (greyc/norns,
521 <https://hub.docker.com/r/greyc/norns>).

522 Pipeline Pilot (BIOVIA Pipeline Pilot, Release 7.5, San Diego: Dassault Systems) is
523 commercial software.

524 **FUNDING**

525 Etienne Lehembre and Johanna Giovannini received funding from ANR. This work was
526 supported by the ANR Involvd project (ANR-20-CE-23-0023).

527 **Author Information**

528 Authors and affiliations

529 Centre d'Etudes et de Recherche sur le Médicament de Normandie, Normandie Université,
530 UNICAEN, CERMN, 14000 Caen, France.

531 Johanna Giovannini (johanna.giovannini@unicaen.fr), Damien Geslin
532 (damien.geslin@unicaen.fr), Alban LEPAILLEUR (alban.lepailleur@unicaen.fr), Ronan
533 Bureau (ronan.bureau@unicaen.fr).

534 Groupe de Recherche en Informatique, Image, Automatique et Instrumentation de Caen,
535 Normandie Université, UNICAEN, ENSICAEN, CNRS, GREYC, 14000 Caen, France

536 Etienne Lehembre (etienne.lehembre@unicaen.fr), Jean-Luc Lamotte (jean-luc.lamotte@unicaen.fr),
537 Abdelkader Ouali (abdelkader.ouali@unicaen.fr), Bertrand Cuissart
538 (bertrand.cuissart@unicaen.fr), Bruno Cremilleux (bruno.cremilleux@unicaen.fr), Albrecht
539 Zimmermann (albrecht.zimmermann@unicaen.fr)

540 Univ. Bordeaux, CNRS, Bordeaux INP, INRIA, LaBRI, Talence, France.

541 David Auber (david.auber@u-bordeaux.fr).

542 Contributions

543 RB, EL and AZ wrote the initial draft. RB, BCuissart, and AZ designed the project (ANR
544 Involvd). EL, RB, BCuissart, AO, and AZ developed the method. EL, BCremilleux, and AO
545 collaborated on program design in accordance with the method's specifications. JLL contributed
546 to the integration of Norns into the program. JG and AL conducted the analysis of BCR-ABL
547 data related to PADs. DG and DA were involved in defining the layouts. All authors have
548 reviewed and approved the manuscript.

549 Corresponding Author.

550 Correspondence to Ronan Bureau

551 **Ethics declarations**

552 Competing interests

553 The authors declare no competing financial interests.

554 **References**

- 555 1. The Practice of Medicinal Chemistry - 4th Edition. [https://www.elsevier.com/books/the-](https://www.elsevier.com/books/the-practice-of-medicinal-chemistry/wermuth/978-0-12-417205-0)
556 [practice-of-medicinal-chemistry/wermuth/978-0-12-417205-0](https://www.elsevier.com/books/the-practice-of-medicinal-chemistry/wermuth/978-0-12-417205-0). Accessed 12 Apr 2023
- 557 2. Langer T, Hoffmann RD (2006) Pharmacophores and Pharmacophore Searches. John
558 Wiley & Sons
- 559 3. Métivier J-P, Cuissart B, Bureau R, Lepailleur A (2018) The Pharmacophore Network: A
560 Computational Method for Exploring Structure-Activity Relationships from a Large
561 Chemical Data Set. *J Med Chem* 61:3551–3564.
562 <https://doi.org/10.1021/acs.jmedchem.7b01890>
- 563 4. Wermuth CG, Ganellin CR, Lindberg P, Mitscher LA (1998) Glossary of terms used in
564 medicinal chemistry (IUPAC Recommendations 1998). *Pure Appl Chem* 70:1129–1143.
565 <https://doi.org/10.1351/pac199870051129>
- 566 5. Lin S-K (2000) Pharmacophore Perception, Development and Use in Drug Design. Edited
567 by Osman F. Güner. *Molecules* 5:987–989. <https://doi.org/10.3390/50700987>
- 568 6. Daveu C, Bureau R, Baglin I, et al (1999) Definition of a Pharmacophore for Partial
569 Agonists of Serotonin 5-HT₃ Receptors. *J Chem Inf Comput Sci* 39:362–369.
570 <https://doi.org/10.1021/ci980153u>
- 571 7. Scior T, Bernard P, Medina-Franco JL, Maggiora GM (2007) Large compound databases
572 for structure-activity relationships studies in drug discovery. *Mini Rev Med Chem* 7:851–
573 860. <https://doi.org/10.2174/138955707781387858>
- 574 8. Horvath D (2008) Chapter 2: Topological Pharmacophores. In: *Chemoinformatics*
575 *Approaches to Virtual Screening*. pp 44–75
- 576 9. Schneider G, Neidhart W, Giller T, Schmid G (1999) “Scaffold-Hopping” by Topological
577 Pharmacophore Search: A Contribution to Virtual Screening. *Angew Chem Int Ed*
578 38:2894–2896. [https://doi.org/10.1002/\(SICI\)1521-3773\(19991004\)38:19<2894::AID-](https://doi.org/10.1002/(SICI)1521-3773(19991004)38:19<2894::AID-ANIE2894>3.0.CO;2-F)
579 [ANIE2894>3.0.CO;2-F](https://doi.org/10.1002/(SICI)1521-3773(19991004)38:19<2894::AID-ANIE2894>3.0.CO;2-F)
- 580 10. Cheng H, Yan X, Han J, Hsu C-W (2007) Discriminative Frequent Pattern Analysis for
581 Effective Classification. In: 2007 IEEE 23rd International Conference on Data
582 Engineering. pp 716–725
- 583 11. Blumenthal DB, Boria N, Gamper J, et al (2020) Comparing heuristics for graph edit
584 distance computation. *VLDB J* 29:419–458. <https://doi.org/10.1007/s00778-019-00544-1>
- 585 12. Geslin D, Lepailleur A, Manguin J-L, et al (2022) Deciphering a Pharmacophore Network:
586 A Case Study Using BCR-ABL Data. *J Chem Inf Model* 62:678–691.
587 <https://doi.org/10.1021/acs.jcim.1c00427>
- 588 13. Hasse Diagram - an overview | ScienceDirect Topics.
589 <https://www.sciencedirect.com/topics/mathematics/hasse-diagram>. Accessed 13 Jun 2023

- 590 14. Davey BA, Priestley HA (2002) Introduction to Lattices and Order, 2nd ed. Cambridge
591 University Press, Cambridge
- 592 15. Ganter B, Wille R (1999) Concept Lattices of Contexts. In: Ganter B, Wille R (eds) Formal
593 Concept Analysis: Mathematical Foundations. Springer, Berlin, Heidelberg, pp 17–61
- 594 16. Gaulton A, Hersey A, Nowotka M, et al (2017) The ChEMBL database in 2017. *Nucleic
595 Acids Res* 45:D945–D954. <https://doi.org/10.1093/nar/gkw1074>
- 596 17. Bento AP, Gaulton A, Hersey A, et al (2014) The ChEMBL bioactivity database: an
597 update. *Nucleic Acids Res* 42:D1083–D1090. <https://doi.org/10.1093/nar/gkt1031>
- 598 18. BCR/ABL - BCR/ABL protein - Homo sapiens (Human) | Publications | UniProtKB |
599 UniProt. <https://www.uniprot.org/uniprotkb/Q16189/publications>. Accessed 27 Oct 2023
- 600 19. ChemBL download version. [https://chembl.gitbook.io/chembl-interface-
601 documentation/downloads](https://chembl.gitbook.io/chembl-interface-documentation/downloads). Accessed 26 Oct 2023
- 602 20. Lehembre E, Bureau R, Cremilleux B, et al (2022) Selecting Outstanding Patterns Based
603 on Their Neighbourhood. In: Bouadi T, Fromont E, Hüllermeier E (eds) *Advances in
604 Intelligent Data Analysis XX*. Springer International Publishing, Cham, pp 185–198
- 605 21. Fournier-Viger P, Gueniche T, Zida S, Tseng VS (2014) ERMiner: Sequential Rule
606 Mining Using Equivalence Classes. In: Blockeel H, van Leeuwen M, Vinciotti V (eds)
607 *Advances in Intelligent Data Analysis XIII*. Springer International Publishing, Cham, pp
608 108–119
- 609 22. Xing L, Klug-Mcleod J, Rai B, Lunney EA (2015) Kinase hinge binding scaffolds and
610 their hydrogen bond patterns. *Bioorg Med Chem* 23:6520–6527.
611 <https://doi.org/10.1016/j.bmc.2015.08.006>
- 612 23. Dogrusoz U, Giral E, Cetintas A, et al (2009) A layout algorithm for undirected compound
613 graphs. *Inf Sci* 179:980–994. <https://doi.org/10.1016/j.ins.2008.11.017>
- 614 24. Shannon P, Markiel A, Ozier O, et al (2003) Cytoscape: a software environment for
615 integrated models of biomolecular interaction networks. *Genome Res* 13:2498–2504.
616 <https://doi.org/10.1101/gr.1239303>
- 617 25. Massaro F, Molica M, Breccia M (2018) Ponatinib: A Review of Efficacy and Safety.
618 *Curr Cancer Drug Targets* 18:847–856.
619 <https://doi.org/10.2174/1568009617666171002142659>
- 620 26. Ostendorf BN, le Coutre P, Kim TD, Quintás-Cardama A (2014) Nilotinib. *Recent Results
621 Cancer Res Fortschritte Krebsforsch Progres Dans Rech Sur Cancer* 201:67–80.
622 https://doi.org/10.1007/978-3-642-54490-3_3
- 623 27. Waller CF (2018) Imatinib Mesylate. *Recent Results Cancer Res Fortschritte Krebsforsch
624 Progres Dans Rech Sur Cancer* 212:1–27. https://doi.org/10.1007/978-3-319-91439-8_1

625 28. sklearn.model_selection.KFold, version 1.3.2. In: Scikit-Learn. [https://scikit-](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.KFold.html)
626 [learn.org/stable/modules/generated/sklearn.model_selection.KFold.html](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.KFold.html). Accessed 11
627 Apr 2023

628