

This is the peer reviewed version of the following article:

Chemometric-assisted cocrystallization: Supervised pattern recognition for predicting the formation of new functional cocrystals / Fornari, Fabio; Montisci, Fabio; Bianchi, Federica; Cocchi, Marina; Carraro, Claudia; Cavaliere, Francesca; Cozzini, Pietro; Peccati, Francesca; Mazzeo, Paolo P.; Riboni, Nicolò; Careri, Maria; Bacchi, Alessia. - In: CHEMOMETRICS AND INTELLIGENT LABORATORY SYSTEMS. - ISSN 0169-7439. - 226:(2022), pp. 104580-104612. [10.1016/j.chemolab.2022.104580]

*Terms of use:*

The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

14/05/2024 03:57

(Article begins on next page)

1 **Chemometric-Assisted Cocrystallization: Supervised Pattern Recognition for Predicting the**  
2 **Formation of New Functional Cocrystals**

3

4 Fabio Fornari <sup>a</sup>, Fabio Montisci <sup>a</sup>, Federica Bianchi <sup>a,b,\*</sup>, Marina Cocchi <sup>c</sup>, Claudia Carraro <sup>a</sup>, Francesca Cavaliere <sup>d</sup>, Pietro  
5 Cozzini <sup>d</sup>, Francesca Peccati <sup>e</sup>, Paolo P. Mazzeo <sup>a,f</sup>, Nicolò Riboni <sup>a</sup>, Maria Careri <sup>a,g</sup>, Alessia Bacchi <sup>a,f</sup>

6

7 <sup>a</sup> *University of Parma, Department of Chemistry, Life Sciences and Environmental Sustainability, Parco Area delle Scienze*  
8 *17/A, 43124, Parma, Italy*

9 <sup>b</sup> *University of Parma, Interdepartmental Center for Packaging (CIPACK), Parco Area delle Scienze, 43124, Parma, Italy*

10 <sup>c</sup> *University of Modena and Reggio Emilia, Department of Chemical and Geological Sciences, Via Giuseppe Campi 103,*  
11 *41125, Modena, Italy*

12 <sup>d</sup> *University of Parma, Department of Food and Drug, Parco Area delle Scienze 17/A, 43124, Parma, Italy*

13 <sup>e</sup> *Basque Research and Technology Alliance (BRTA), Center for Cooperative Research in Biosciences (CIC bioGUNE),*  
14 *Bizkaia Technology Park 801A, 48160, Derio, Spain*

15 <sup>f</sup> *University of Parma, Biopharmanet-TEC, Parco Area delle Scienze 27/A, 43124, Parma, Italy*

16 <sup>g</sup> *University of Parma, Interdepartmental Center on Safety, Technologies, and Agri-Food Innovation (SITEIA.PARMA), Parco*  
17 *Area delle Scienze, 43124, Parma, Italy*

18

19

20

21

22

23

24 \* Author to whom correspondence should be addressed:

25 e-mail address: federica.bianchi@unipr.it; University of Parma, Department of Chemistry, Life Sciences and

26 Environmental Sustainability, Parco Area delle Scienze 17/A, 43124, Parma, Italy

## 27 **Abstract**

28 Owing to the antimicrobial and insecticide properties, the use of natural compounds like essential oils and their active  
29 components has proven to be an effective alternative to synthetic chemicals in different fields ranging- from drug delivery to  
30 agriculture and from nutrition to food preservation. Their limited application due to the high volatility and scarce water  
31 solubility can be expanded by using crystal engineering approaches to tune some properties of the active molecule by  
32 combining it with a suitable partner molecule (coformer). However, the selection of coformers and the experimental effort  
33 required for discovering cocrystals are the bottleneck of cocrystal engineering. This study explores the use of chemometrics  
34 to aid the discovery of cocrystals of active ingredients suitable for various applications. Partial Least Squares–Discriminant  
35 Analysis is used to discern cocrystals from binary mixtures based on the molecular features of the coformers. For the first  
36 time, by including failed cocrystallization data and considering a variety of chemically diverse compounds, the proposed  
37 method resulted in a successful prediction rate of 83.85% for the test set in the model validation phase and of 62.74% for the  
38 external test set.

39

40

41

42

43

44

45

46

47

48

49

## 50 **Keywords**

51 cocrystal, crystal engineering, chemoinformatics, chemometrics, partial least square discriminant analysis, Quantitative  
52 Structure–Property Relationship

53 **List of abbreviations (sorted alphabetically)**

54 ACC<sub>%</sub>: Accuracy

55 BM: Binary Mixture

56 CC: Cocrystal

57 CCDC: Cambridge Crystallographic Data Centre

58 CSD: Cambridge Structural Database

59 EO: Essential Oil

60 FDA: Food and Drug Administration

61 GRAS: Generally Recognized As Safe

62 LV: Latent Variable

63 MEP: Molecular Electrostatic Potential

64 ~~NER<sub>%</sub>: Non-Error classification Rate~~

65 PC: Principal Component

66 PCA: Principal Component Analysis

67 PLS-DA: Partial Least Squares-Discriminant Analysis

68 PXRD: Powder X-Ray Diffraction

69 QSPR: Quantitative Structure-Property Relationship

70 SEN<sub>%</sub>: Sensitivity

71 VIP: Variable Importance in Projection.

72

73

74

75

76

77

78

79

80

81

## 82 **1. Introduction**

83 In the last few decades, the use of agrochemicals and food preservatives has grown exponentially as a direct consequence of  
84 the rapid increase of the world population [1,2]. Owing to their potential adverse effects on both human health and  
85 environment [2–4], alternative strategies based on the use of more sustainable chemicals have been proposed to support the  
86 food system. Being able to exert antimicrobial, insecticidal, and antioxidant properties [5,6], essential oils (EOs) and their  
87 active components have been used as green substitutes of synthetic chemicals to extend the shelf-life of foodstuff and in pests  
88 control [7–9]. These compounds are Generally Recognized As Safe (GRAS) by the Food and Drug Administration (FDA)  
89 [10], however, despite their appealing properties, their use is limited by their high volatility and poor stability [7,9,11,12].

90 In fact, physicochemical properties of materials play a key role in determining whether a chemical is suitable for a specific  
91 purpose, thus strongly affecting its field of application. Scientists have always desired to obtain materials with target  
92 properties, and crystal engineering is one of the most interesting approaches to synthesize a great variety of crystalline  
93 materials for applications in various fields, ranging from pharmaceuticals to agrochemicals, and from nutraceuticals to  
94 cosmetics [13–16]. The basic idea of crystal engineering is related to the possibility of controlling the crystal structure of  
95 molecules and, therefore, the properties of the resulting solids. Polymorphism, vitrification and cocrystallization are some of  
96 the available strategies to modify the intrinsic properties of molecules without the need of synthetic modifications [17–21].

97 Cocrystals are multicomponent crystalline solid materials in which the constituents (i.e., coformers) are bound in a well-  
98 defined stoichiometric ratio [22,23] *via* non-covalent interactions (e.g., hydrogen bonds, halogen bonds,  $\pi$ - $\pi$  stacking) within  
99 the same crystal structure. Cocrystallization allows for the combination of the desired molecule of interest with properly  
100 selected partner molecules, paving the way to an array of potential materials with enhanced properties [19,24,25]. Within this  
101 frame of reference, cocrystals based on the active components of EOs have been proposed as active ingredients for food  
102 packaging, agrochemical and pharmaceutical applications [7,19,26,27].

103 Despite the great advantages offered by cocrystallization, the proper degree of complementarity between the two partner  
104 molecules required to obtain crystalline materials with the desired properties is not easy to assess [28–31]. In this context, the  
105 selection of coformers and the great effort required for both the systematic experimental screening and careful characterization  
106 of the products derived from the combination of all the possible coformer pairs represent the major bottleneck of cocrystal  
107 engineering. Computational techniques represent a powerful tool to reduce the experimental effort required for the discovery  
108 of new cocrystals, enabling to evaluate beforehand whether a cocrystal can be obtained starting from pre-selected coformers.

109 These *in silico* strategies can be based on the calculation of a variety of parameters useful to predict the formation of a  
110 cocrystal, such as lattice energy [32], solubility [33], hydrogen bond propensity along with the quantitation of molecular  
111 interaction energy [29,30], and molecular complementarity [34].

112 Despite the massive efforts spent to develop a method to predict cocrystal formation, at present none of the proposed strategies  
113 has proven to be both totally reliable and easy to apply.

114 Chemometrics could play a pivotal role in cocrystal discovery: up to now only few Machine Learning methods have been  
115 proposed in predicting cocrystal formation, enabling the screening of new cocrystals once a supervised model is properly  
116 trained and validated. In the study proposed by Devogelaer et al., information of successful cocrystallization experiments was  
117 directly taken from the Cambridge Structural Database (CSD) [35] and Artificial Neural Networks (ANN) were used to predict  
118 the formation of new cocrystals [36]. Similarly, Wang et al. relied on a consensus method based on multiple Random Forest  
119 algorithms, in which the successful cocrystallization dataset was integrated with randomly generated failed cocrystallization  
120 data [37]. These approaches are reported in the literature as network-based methods. Additional studies were based on the use  
121 of successful and unsuccessful cocrystallization datasets obtained from experimentation, literature, and/or the CSD for  
122 screening specific classes of cofomers. Within this framework, Przybyłek et al. used Multivariate Adaptive Regression  
123 Splines to predict the formation of dicarboxylic and phenolic acid-based cocrystals [38,39], whereas Wicker et al. focused on  
124 variously substituted benzoic acids and benzamides using a Support Vector Machine algorithm [24]. Vriza et al. used an  
125 ensemble one-class classification method to aid the discovery of  $\pi$ - $\pi$  cocrystals, thus giving a great contribution in enriching  
126 one of the most under-represented classes of cocrystals in the CSD [40]. Most recently, Mswahili et al. developed a cocrystal  
127 screening method based on ANN by using both successful and unsuccessful experimental cocrystallization data retrieved from  
128 the literature and a plethora of molecular descriptors calculated using Mordred [41,42].

129 In the frame of a research activity dealing with the synthesis of new functional cocrystals based on the active constituents of  
130 EOs and other GRAS molecules to broaden their applicability in the industrial field [7,19], we propose a chemometric  
131 approach to aid the discovery of new cocrystalline materials.

132 For the first time, a training set based on the results of failed (binary mixtures, BM) and successful (cocrystal, CC)  
133 cocrystallization experiments was used for the computation of a Quantitative Structure–Property Relationship-like (QSPR)  
134 model based on Partial Least Squares–Discriminant Analysis (PLS–DA), after preliminary exploratory analysis by Principal  
135 Component Analysis (PCA). The PLS-DA approach, with respect to network-based methods offers the advantages of having

136 only one parameter to optimize, i.e., the model dimensionality, and the direct interpretation of the importance of descriptors  
137 in classification, while highlighting their interplay (by inspection of weights and loadings plots).  
138 The effectiveness of the study relies on the use of compounds belonging to different chemical classes and a reduced number  
139 of 1D, 2D, and 3D molecular descriptors of various nature (e.g., constitutional, geometric, physical, topological, and surface  
140 area-based descriptors) [24,38,39,43], enabling the high-throughput screening of novel cocrystalline materials and offering  
141 maximum flexibility and effectiveness at a minimum computational and experimental cost.

142

## 143 **2. Experimental Procedures**

### 144 2.1. Mechanochemical protocol and class assignation

145 All the molecules in the dataset were chosen among the list of GRAS molecules drawn up by the FDA [10]. Selected pair of  
146 molecules among the chosen ones were assigned either to the CC class or to the BM class. Pairs of molecules in the dataset  
147 for which a cocrystal structure was already described in literature were individuated in the Cambridge Structural Database  
148 (CSD) [35] with the Cambridge Crystallographic Data Centre (CCDC) software ConQuest [44] and visualized with Mercury  
149 [45]. They are reported in Section 3 of the Supplementary Material, together with their unique CSD refcode and the reference  
150 to the original publications.

151 Cocrystallization for all the pairs with no known structure in literature was instead attempted with the following  
152 mechanochemical protocol. All the reagents employed were commercially available and used as such in all the experiments.  
153 Equimolar amounts of each reagent were directly mixed in an agate mortar and subjected to manual grinding for 10-15  
154 minutes, without using any solvent. The resulting powder samples were collected in closed vials. Assignment to CC or BM  
155 classes was performed by comparing the Powder X-ray Diffraction (PXRD) pattern of the ground sample with those of the  
156 pure reagents. Possible occurrence of polymorphic transitions for the reagents was excluded by comparing the experimental  
157 PXRD data after milling with the calculated pattern of all the known crystalline forms of the reagents. The occurrence of new  
158 peaks, unexplained by the presence of unreacted reagents, was taken as indication that cocrystallization had occurred and the  
159 sample was assigned to the CC class. In case no additional peaks appeared in the PXRD pattern the sample was instead  
160 classified as a BM.

161

162

## 163 2.2. Powder X-ray diffraction

164 Typically, PXRD data were collected on a Rigaku Smartlab XE diffractometer in  $\theta$ - $\theta$  Bragg-Brentano geometry with Cu K $\alpha$   
165 radiation. The samples were placed on glass supports and exposed to radiation ( $1.5^\circ \leq 2\theta \leq 50^\circ$ ) at a scan rate of  $10^\circ/\text{min}$ . The  
166 diffracted beam was collected on a 2D Hypix 3000 solid state detector.  $5^\circ$  radiant soller were used as a compromise for high  
167 flux and moderate peak asymmetry at low angles. Beam stopper and anti-scatterer air component were used to mitigate the  
168 profile at low angle. In some rare cases, the data were collected on a Thermo Fisher Scientific ARL X'TRA diffractometer in  
169  $\theta$ - $\theta$  Bragg-Brentano geometry with Cu K $\alpha$  radiation ( $3^\circ \leq 2\theta \leq 30^\circ$  at a scan rate of  $5^\circ/\text{min}$ , or  $3^\circ \leq 2\theta \leq 40^\circ$  at a scan rate of  
170  $0.3^\circ/\text{min}$ ).

171

## 172 3. Computational Methods

### 173 3.1. Molecular descriptors calculation

174 For each molecule 31 molecular descriptors were calculated (Table S1). A theoretical background for the less known  
175 descriptors is given in Section 4 of the Supplementary Material.

176 The molecular weight, the number of atoms, the number of bonds, the number of hydrogen bond donor sites and the number  
177 of hydrogen bond acceptor sites were calculated with FLAP software (Fingerprint for Ligand and Protein) [46] at pH 7.0,  
178 using the 3D structures of all molecules in SDF format as input (downloaded from the PubChem database). The number of  
179 rotatable bonds, the number of rings, the hydrophobicity (accounted as the number of hydrophobic centers), the logP  
180 (logarithm of octanol/water partition coefficient), the molecular volume, the total molecular dipole moment (based on point  
181 charge distribution in the molecule), and its components along the axes (using the principal axes of the molecular graph) were  
182 then calculated for the same structures using Sybyl 8.1 [47] (www.tripos.com) and taking in consideration the protonation  
183 state of molecules. The same software was also used to estimate the strain energy of the molecule without performing any  
184 geometry optimization. This energy term relies on an electrostatic calculation from atomic charges using the internal Tripos  
185 force field [48]. For the estimation of molecular volume and dipole moment, a specific SPL script was employed. The  
186 calculated volume is enclosed in a water-accessible surface computed at a repulsive interaction energy of 0.20 kcal/mol with  
187 a water probe. A custom Python script was used to automatically calculate the Solvent Accessible Surface Area (SASA) in  
188 PyMol 2.0 [49], with the dot density parameter set to 4. The number of heteroatoms, the number of valence electrons, and the  
189 indexes  ${}^0\chi$ ,  ${}^0\chi^n$ ,  ${}^0\chi^v$ ,  $\alpha_{HK}$ ,  ${}^1\kappa_a$ , LabuteASA, SMR\_VSA, PEOE\_VSA, and TPSA were obtained running a Python 3.7 code with



190 the open-source cheminformatics toolkit RDKit Q4 2013 [50]. The average isotropic polarizability  $\alpha_{iso}$ , the polarizability  
191 anisotropy  $\Delta\alpha$ , and the Molecular Electrostatic Potential (MEP) were calculated with Gaussian 16 [51] following the *in vacuo*  
192 Density-Functional Theory optimization of all the molecules, employing the hybrid functional B3LYP and the People double-  
193 z basis set 6-31+g(d,p).

194 Postprocessing of the MEP to extract critical points at a given electron density isosurface was performed with a custom Python  
195 3.6.1 script on a three-dimensional map (cube format) with a sampling density of 6 points/Bohr along the three directions.  
196 The MEP was analyzed at an electron density isosurface of 0.002 a.u. with a tolerance of 0.001 a.u., meaning that only MEP  
197 values corresponding to regions of space with electron density in the 0.001–0.003 a.u. range were considered. A first set of  
198 critical points was identified comparing MEP values of each cube point with those of its 6 nearest neighbors. A point was  
199 considered a local minimum if the number of nearest neighbors with higher MEP was greater or equal to a given integer (4).  
200 Likewise, a point was considered a local maximum if the number of nearest neighbors with lower MEP was greater or equal  
201 to the same integer. This first step yielded a large number of candidate critical points encompassing a wide range of MEP  
202 values. Since our focus was on identifying the regions of the molecules likely to be involved in strong hydrogen bonds within  
203 the cocrystal, in a second step this first set of points was filtered based on the MEP values of the global minimum and  
204 maximum. This was done as follows: for each local minimum (maximum), the ratio between its MEP value and that of the  
205 global minimum (maximum) was computed, and the point was kept only if the ratio exceeded a given threshold (0.1). In this  
206 way, only points corresponding to shallow critical points were discarded. This step allowed to identify the MEP isosurface  
207 regions corresponding to hydrogen bond donors and acceptors. However, due to the rugged character of the MEP map,  
208 multiple critical points of the same type could appear in close proximity. To univocally map a given region of the isosurface  
209 to a MEP value, critical points close to each other (below a distance threshold of 1.0 Bohr) were merged iteratively, keeping  
210 only the lower MEP point for minima and higher MEP point for maxima. The algorithm then provided the final set of MEP  
211 critical points at the given electron density isosurface.

212

### 213 3.2. Data analysis

214 The entire data analysis was carried out in MATLAB R2019a environment (Mathworks, Natick, Massachusetts, USA) with  
215 the aid of the- PLS\_Toolbox 8.7.1 (Eigenvector Research Inc., Washington, USA) chemometric package.

~~was used to carry out preprocessing, Principal Component Analysis (PCA) and PLS-DA computation, and to split the original dataset into calibration and test set. The proper number of latent variables (LVs) to be retained was evaluated by running a homemade MATLAB routine.~~

### 3.2.1. Data preprocessing

Each sample was described by  $m = 31$  variables (Table S1): the ~~absolute value of the difference between~~ ~~difference in absolute value of~~ the molecular descriptors of the two partner molecules was calculated, thus obtaining the predictor matrix  $X$  ( $181 \times 31$ ). The class membership was binary encoded (1: belonging to the class; 0: otherwise) in a dummy matrix  $Y$  ( $181 \times 2$ ) with each column representing one of the two modelled classes. The dataset was split in two subsets by using the ~~Kennard-Stone~~ ~~samplingDuplex~~ algorithm [52]: ~~8070%~~ of the data were used as calibration set,  $X_{\text{cal}}$  (~~146-127~~  $\times 31$ ) and  $Y_{\text{cal}}$  (~~146-127~~  $\times 2$ ), whereas the remaining ~~2030%~~ were used as test set,  $X_{\text{test}}$  (~~35-54~~  $\times 31$ ) and  $Y_{\text{test}}$  (~~35-54~~  $\times 2$ ). The calibration set and the test set are reported in Table S2 and Table S3, respectively.

Before carrying out both exploratory multivariate data analysis and the computation of the supervised model, the calibration matrix  $X_{\text{cal}}$  was preprocessed column-wise by performing mean centering and scaling to unit variance. Mean centering was applied on the response matrix  $Y_{\text{cal}}$  to ensure the stability of the model.

### 3.2.2. Exploratory multivariate data analysis

PCA [53–55] was carried out preliminarily on the calibration set  $X_{\text{cal}}$  to assess the distribution of the samples and to check for potential data structures. Reduction of data dimensionality is carried out through the linear combination of the original variables in a set of orthogonal ones, i.e., Principal Components (PCs), which identify the direction of maximum variance. This is summarized in the decomposition equation:

$$X_{\text{cal}} = TP^T + E$$

where  $T$  and  $P$  represent, respectively, the coordinates of the samples projected in the reduced space, i.e., the scores, and the weights each original variable has on a given PC, i.e., the loadings. The deviations from the model are accounted in the error matrix  $E$ .

### 243 3.2.3. Supervised pattern recognition

244 PLS–DA [56,57] was used to discriminate pairs of partner molecules whose combination forms CCs from the ones giving  
245 BMs. PLS–DA is based on PLS regression [58]. Briefly, this supervised technique decomposes the predictor matrix  $X_{\text{cal}}$  and  
246 the dependent variables matrix  $Y_{\text{cal}}$  in a PCA-like way and imposes inner linear relationships between the  $X_{\text{cal}}$  and the  $Y_{\text{cal}}$   
247 scores as follows:

$$248 \quad \quad \quad U = bT$$

249 where  $T$  and  $U$ , are the  $X_{\text{cal}}$  and the  $Y_{\text{cal}}$  scores, respectively. This is accomplished by rotating the Latent Variable (LV) space  
250 of  $X_{\text{cal}}$  through a weight matrix  $W$  in a way that maximizes the covariance between  $T$  and  $U$ . The PLS regression model is  
251 summarized as:

$$252 \quad \quad \quad Y_{\text{cal}} = X_{\text{cal}}B + E$$

253 where  $E$  is the error matrix and  $B$  is the pseudo-regression coefficient matrix expressed according to the following equation:

$$254 \quad \quad \quad B = W(P^T W)^{-1} \text{diag}(b) Q$$

255 where  $P$  and  $Q$  are the  $X_{\text{cal}}$  and the  $Y_{\text{cal}}$  loadings, respectively.

256 In this case, the dependent variables in the  $Y_{\text{cal}}$  matrix are defined as dummy variables, one for each modelled class, taking  
257 values of 1 if the sample belongs to the class and 0 otherwise. Current implementation of PLS–DA may differ on the basis of  
258 how the classification rule is defined. In this work, a pure discriminant rule (samples are assigned univocally to only one  
259 category) was applied, and thus a sample is assigned to the class for which the predicted response  $\hat{y}$  is the highest (i.e.,  $\hat{Y}$   
260 values are continuous and not dummy as they were codified).

261 The proper number of LVs was chosen according to the maximum Non-Error classification Rate accuracy (NERACC%; i.e.,  
262 the percentage of samples correctly assigned to the respective class) in leave-more-out cross validation, adopting a Venetian  
263 blinds cancellation scheme with 10 splits (blind thickness: 1). This operation was carried out by running a custom MATLAB  
264 routine. The performance of the classification model was evaluated both on the calibration and the test sets in terms of  
265 NERACC% as well as showing the confusion matrix. In addition, the sensitivity (SEN%, i.e., the percentage of samples within  
266 a class that were correctly assigned to their class) was calculated for both classes.

267 The importance of each predictor was estimated in terms of Variable Importance in Projection (VIP score) [56]. The VIP  
268 score of the  $j^{\text{th}}$  variable in the  $X$  space is defined as the component-wise sum of its PLS weight  $w_j$  on the  $f^{\text{th}}$  component  
269 multiplied by the fraction of variance of the  $Y$  explained by that component, according to the following equation:

$$VIP_j^2 = \frac{1}{SSY_{\text{tot}}^F} \sum_{f=1}^F w_{jf}^2 SSY_f J$$

where  $J$  is the number of variables in the  $X$  space and  $F$  is the number of LVs that were retained. Since:

$$\sum_{j=1}^J VIP_j^2 = J$$

the proposed threshold for determining whether a variable could be considered important is set to 1..

Finally, the predictive capability of the model was evaluated on an external set of  $N = 58$  binary combinations of partner molecules. An overview of the involved samples is reported in Table S4 along with their estimated  $\hat{y}$  values.

The number of entries in the BM and CC classes for training, test, and external validation sets is reported in Table 1.

**Table 1.** Number of CC and BM samples in the calibration, test, and external validation set. The last column reports the total number  $N$  of samples *per set*.

	<u>CC</u>	<u>BM</u>	<u><math>N</math></u>
<u>Calibration set</u>	<u>71</u>	<u>56</u>	<u>127</u>
<u>Test set</u>	<u>30</u>	<u>24</u>	<u>54</u>
<u>External validation set</u>	<u>31</u>	<u>27</u>	<u>58</u>

## 4. Results and Discussion

The cocrystallization experiments were carried out mechanochemically by manual neat grinding of the two substances. This method was selected among many possible others due to its simplicity and promptness, allowing us to screen several molecular pairs in a standardized way [21,59]. The classification as BM or CC was assessed by PXRD patterns (available in Section 3 of the Supplementary Material). Cocrystals already present in the CSD were also included in our dataset.

### 4.1. Data analysis

Data handling prior to analysis can affect the way the model is trained, thus having consequences on its interpretation.

Since samples are the result of the combination of two partner molecules, each one described by its own set of descriptors, the concatenation strategy, i.e., listing the descriptors of the first coformer followed by those of the second coformer, was

291 discarded due to the lack of commutation between the two sets of descriptors. In fact, the order in which the molecular  
292 descriptors are listed represents an *a priori* decision on which compound is acting as molecule of interest or partner molecule.  
293 In our case this would be sub-optimal since many of the molecules in our dataset could assume both roles. Therefore, in order  
294 to address the problem described above a commutative strategy capable of avoiding the production of indeterminate forms  
295 should be chosen. Considering that in the dataset many constitutional molecular descriptors were characterized by a few non-  
296 zero values, the calculation of both products and ratios between the descriptors was discarded since additional zero values or  
297 indeterminate forms could be generated. Also, the division is non-commutative.

298 In order to combine the two partner molecules without imposition on their role, the absolute value of the difference difference  
299 in absolute value between the molecular descriptors of the partner molecules was calculated and used to describe each binary  
300 combination [38,39], giving maximum flexibility to the model. The information achieved is still relevant and easily  
301 interpretable since it is related to the dissimilarity between the descriptors. In fact, the differences in absolute value absolute  
302 value of the differences between descriptors for each case are the elements of the Manhattan distance [60], one of the possible  
303 indexes to account for the dissimilarity between cases in a multivariate way.

304 In the present study, the use of basic linear modelling methods was preferred with respect to non-linear modelling, such as  
305 ANN [61,62], to keep the calculations as simple as possible, and to ensure a certain degree of interpretability of the results.  
306 Furthermore, the number of samples is too limited to ensure proper tuning of the ANN hyperparameters.

307 For a fruitful discussion, samples and variables are reported in the text according to their identification number as follows: i)  
308 samples are written in plain text; ii) variables are underlined. The key is available in Tables S1–S3.

309

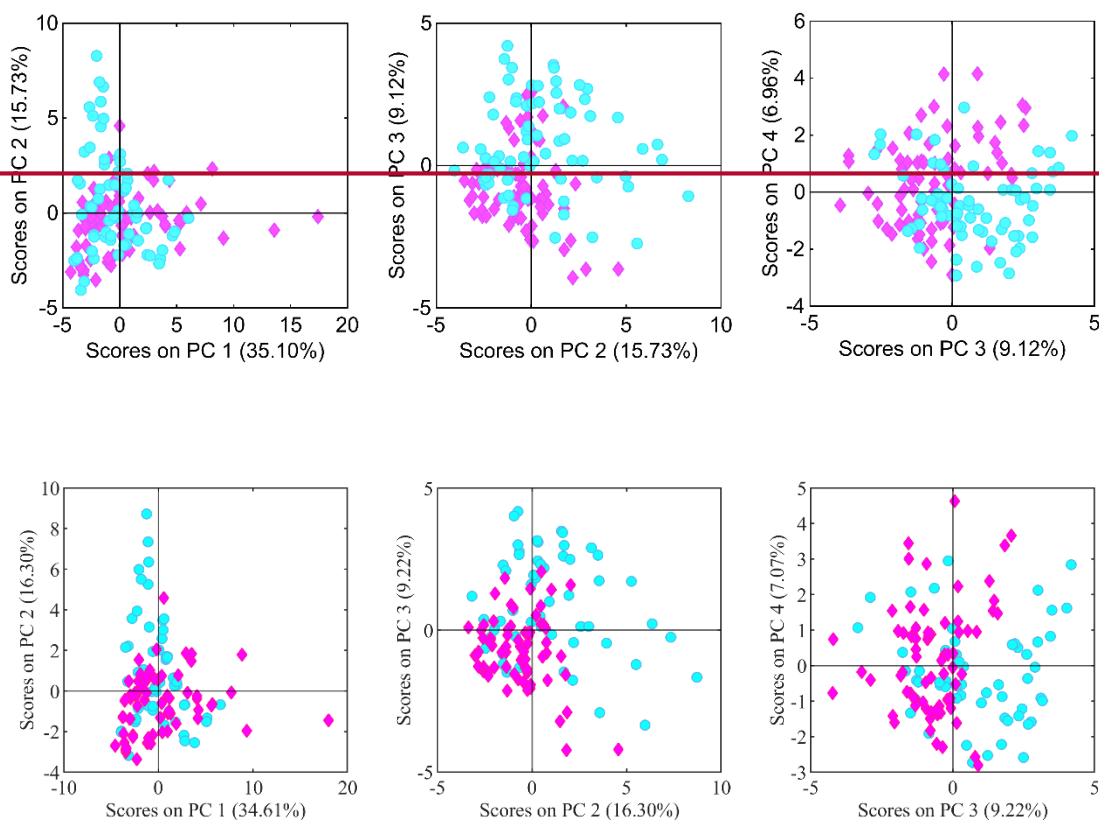
#### 310 *4.1.1. Exploratory multivariate data analysis*

311 After data preprocessing, PCA was used in an exploratory way to assess the presence of potential data structures in the  
312 calibration set.

313 Four PCs were retained explaining 67% of the variance. As shown in the score plots (Figure 1), a mild segregation-separation  
314 was present in the PC 3-2 vs. PC 4-3 score plot, with the groups separated by the bisecting line of the II and the III-IV quadrant.

315 In addition, most of the CC samples occupied the III quadrant of the PC 2 vs. PC 3 score plot and were, in general, less  
316 scattered than BM samples.

317



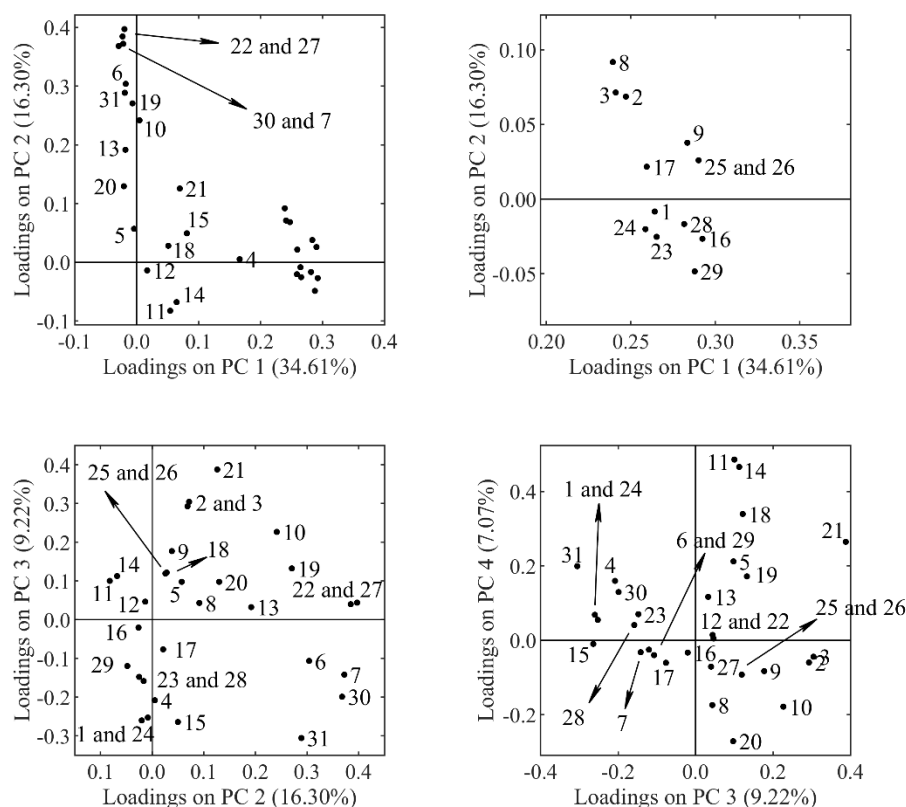
318

319

320 **Figure 1.** Scores of the samples on the first 4 PCs: PC 1 vs. PC 2 (left), PC 2 vs. PC 3 (center), PC 3 vs. PC 4 (right). The  
 321 fraction of variance explained by a given component is reported as a percentage value in parenthesis on the corresponding  
 322 axis. Samples are marked according to their class (empty-magenta diamonds: cocrystals; black-cyan circles: binary mixtures).

323

324 As for the loading plots depicted in Figure [S12](#), it can be observed that PC 1 explains the features related both to differences  
 325 in molecular dimensions (e.g., [1](#), [2](#), [3](#), [9](#), [16](#), [29](#)) and connectivity ([23](#), [24](#), [25](#), [26](#)). PC 2 considers the dissimilarities in the  
 326 electronic properties (e.g., [6](#), [7](#), [19](#), [22](#), [30](#)) of the two partner molecules as well as the difference in their number of  
 327 heteroatoms ([27](#)). PC 3 accounts for more specific features, such as differences in molecular complexity ([21](#)), molecular  
 328 refractivity ([31](#)), and energy ([15](#)). Lastly, PC 4 considers the dissimilarity in: i) component of the dipole along the  $x$  axis ([11](#)),  
 329 ii) total dipole ([14](#)) and iii) minimum of the MEP surface ([20](#)). Regarding the samples belonging to the CC class, these  
 330 compounds were characterized by partner molecules with a similar behavior in terms of molecular complexity ([21](#)) and  
 331 octanol/water partition coefficient ([10](#)).



332

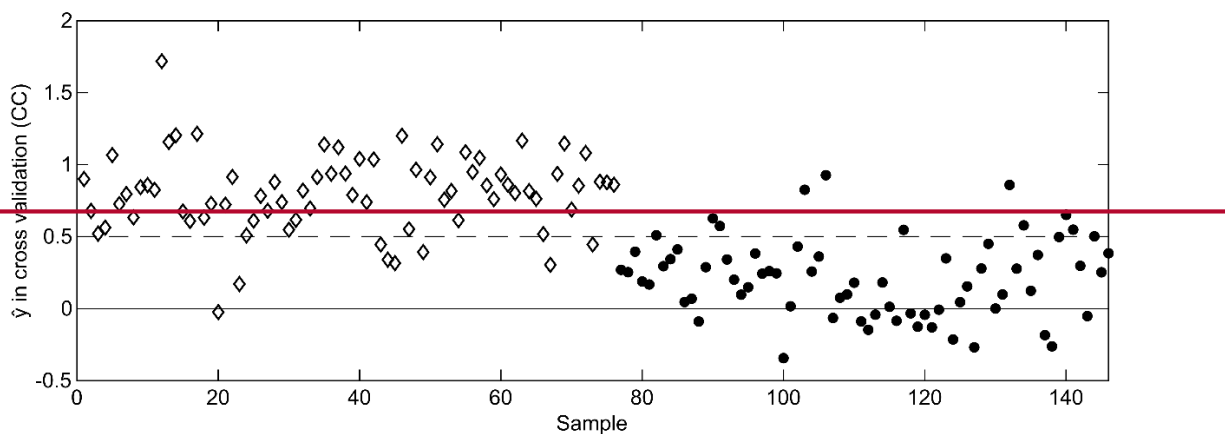
333 **Figure 2.** Loading plots related to the PCA decomposition. PC 1 vs. PC 2 (top-left), magnification of the PC 1 vs. PC 2 (top-  
 334 right), PC 2 vs. PC 3 (bottom-left), PC 3 vs. PC 4 (bottom-right). The fraction of variance explained by a given component is  
 335 reported as a percentage value in parenthesis on the corresponding axis.

336

337 *4.1.2. Supervised pattern recognition*

338 The relationship between the class membership and the variables was exploited by means of PLS-DA. ~~Four~~ Six LVs were  
 339 retained according to maximum NERACC% in cross validation. The PLS-DA model captured the 6372% and 6469% of the  
 340 variance of the  $X_{cal}$  and  $Y_{cal}$ , respectively. The values of the predicted response  $\hat{y}$  in cross validation related to CC samples are  
 341 plotted in Figure 2S1.

342



**Figure 2:**  $\hat{y}$  in cross validation for the CC samples. The dashed horizontal line shows the hard classification threshold of 0.50, halfway between the codified 0 (BM) and 1 (CC). Samples are marked according to their class (empty diamonds: co-crystals; black circles: binary mixtures).

A summary of the performance of the obtained model is reported in the confusion matrix (Table 42), whereas a graphical representation of the estimated and predicted values  $\hat{y}$  for the CC class is reported in Figure S2. As reported in Table 42, all the BM samples belonging to the test set were correctly classified except for one 5 samples, whereas only 53 out of 2530 CC samples were wrongly assigned to the BM class, obtaining a **NERACC%** of 8385%. Similarly, a high **NERACC%** of 92% was obtained when the samples belonging to the calibration set were predicted by the model. The achieved results are extremely satisfactory, allowing for the *a priori* selection of the partner molecules required for the synthesis of novel co-crystals.

**Table 42.** Confusion matrix of the calibration and the test sets for the PLS-DA model. The second-last line shows the SEN% for the modelled classes and last line shows the **NERACC%** for the calibration and the test sets.

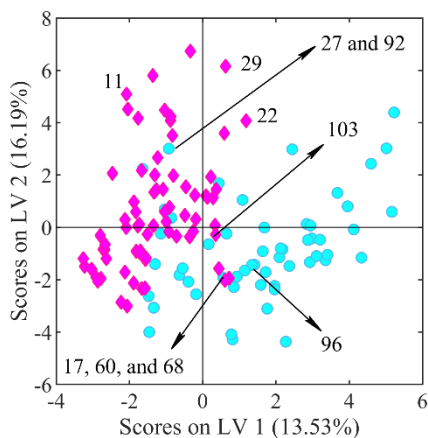
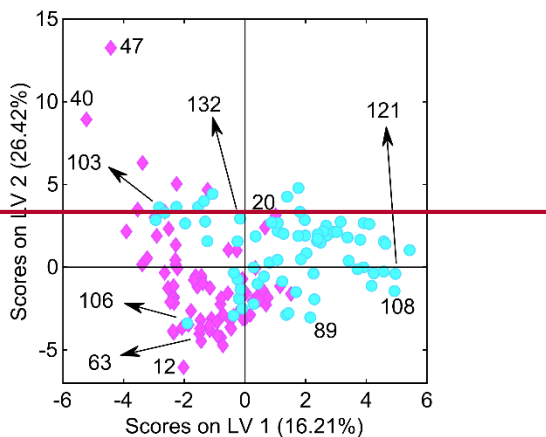
	Calibration set		Test set	
	Predicted as		Predicted as	
	CC	BM	CC	BM
True CC	<del>7166</del>	5	<del>2027</del>	<del>53</del>
True BM	<del>75</del>	<del>6351</del>	<del>45</del>	<del>919</del>
<b>SEN%</b>	<b>93%</b>	<b>91%</b>	<b>90%</b>	<b>79%</b>



358

359 The distribution of the samples in the reduced space of the LVs can be observed by inspecting the score plot (Figure 3),  
360 whereas information regarding suspicious and/or influential samples can be retrieved by the squared residuals  $Q$  vs.  
361 Hotelling's  $T^2$  and the residuals vs. leverage plots (Figure 4).

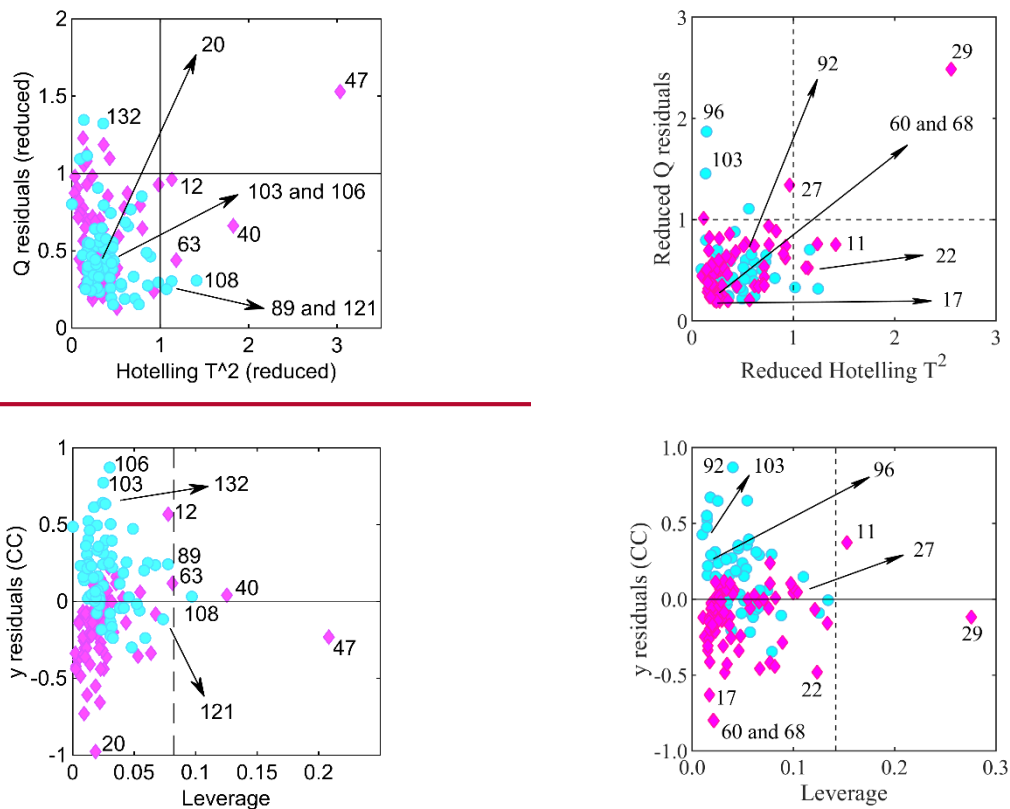
362



364

365 **Figure 3.** Distribution of the samples in the score space (LV 1 vs. LV 2). The fraction of variance explained by a given  
366 component is reported as a percentage value in parenthesis on the corresponding axis. Samples are marked according to their  
367 class (empty magenta diamonds: cocryystals; black-cyan circles: binary mixtures).

368



369

370 **Figure 4.** Reduced squared residuals  $Q$  vs. reduced Hotelling's  $T^2$  plot, the dashed horizontal and vertical lines show the  
 371 amplitude of the 95% confidence interval for both parameters (top). Residuals of CC samples vs. leverage plot, the dashed  
 372 vertical line shows the leverage limit (bottom). Samples are marked according to their class (empty magenta diamonds:  
 373 cocrystals; black cyan circles: binary mixtures).

374

375 The maximum separation in the score space was provided by the first two LVs, with the CC samples well grouped/located  
 376 mainly at negative scores both on LV 1 and LV 2, especially in the III quadrant. By contrast, BM samples were more scattered  
 377 and localized mostly on positive scores on LV 1.

378 ~~A peculiar behavior was observed for samples (40) and (47): sample 40 showed Hotelling's  $T^2$  value outside the 95%~~  
 379 ~~confidence interval, whereas sample 47 showed both high Hotelling's  $T^2$  and high squared residuals  $Q$  together with high~~  
 380 ~~leverage in the  $Y_{cal}$  space. These CC samples were obtained by pairing fatty acids (lauric acid, 40, and palmitic acid, 47) with~~  
 381 ~~low molecular weight coformers, i.e., pyrazine and nicotinamide, respectively. The variables responsible for this behavior~~  
 382 ~~were related to the discrepancy in molecular dimensions between the two partner molecules.~~

383 ~~Another BM sample, i.e., the limonene/ascorbic acid (108), was characterized by high leverage and high Hotelling's  $T^2$  value,~~  
384 ~~due to different behavior in terms of hydrogen bond propensity and octanol/water partition coefficient. A similar behavior~~  
385 ~~was observed also by two additional non-influential BM samples based on ascorbic acid paired with cinnamaldehyde (89)~~  
386 ~~and menthone (121).~~

387 ~~Finally, also the tartaric acid/pyrazine (63) and the adipic acid/hexamethyleneamine (12) CC samples showed Hotelling's  $T^2$~~   
388 ~~values outside the 95% confidence interval and the latter resulted in having the highest residuals among CC samples.~~

389 ~~Nevertheless, all these anomalous samples were correctly assigned in cross validation, and the variables with high contribution~~  
390 ~~(not shown) on their Hotelling's  $T^2$  and high squared residuals  $Q$  were characterized by low PLS weights.~~

391 ~~On the other hand, one not anomalous CC sample and two not anomalous BM samples, namely carveol/isonicotinamide (20),~~  
392 ~~eugenol/pyrazine (103), geraniol/menthol (106), respectively, were characterized by very high residuals in absolute value,~~  
393 ~~thus being misclassified. This behavior can be ascribed to the fact that their features were inversely related to their respective~~  
394 ~~class. Also, the BM sample urea/hexamethylenetetramine (132) was misclassified: its features did not align with those of the~~  
395 ~~other samples; in fact, this sample hold high  $Q$  squared residuals. A peculiar behavior was observed for samples 27 and 29.~~

396 ~~Sample 27 showed squared residual  $Q$  outside the 95% confidence interval, whereas sample 29 showed both high Hotelling's~~  
397  ~~$T^2$  and high squared residuals  $Q$  together with high leverage in the  $Y_{cal}$  space. These CC samples were obtained by pairing~~  
398 ~~fatty acids (lauric acid, 27, and palmitic acid, 29) with a low-molecular weight cofomer, i.e., nicotinamide. The variables~~  
399 ~~responsible for this behavior can be related to the discrepancy in molecular dimensions between the two partner molecules.~~

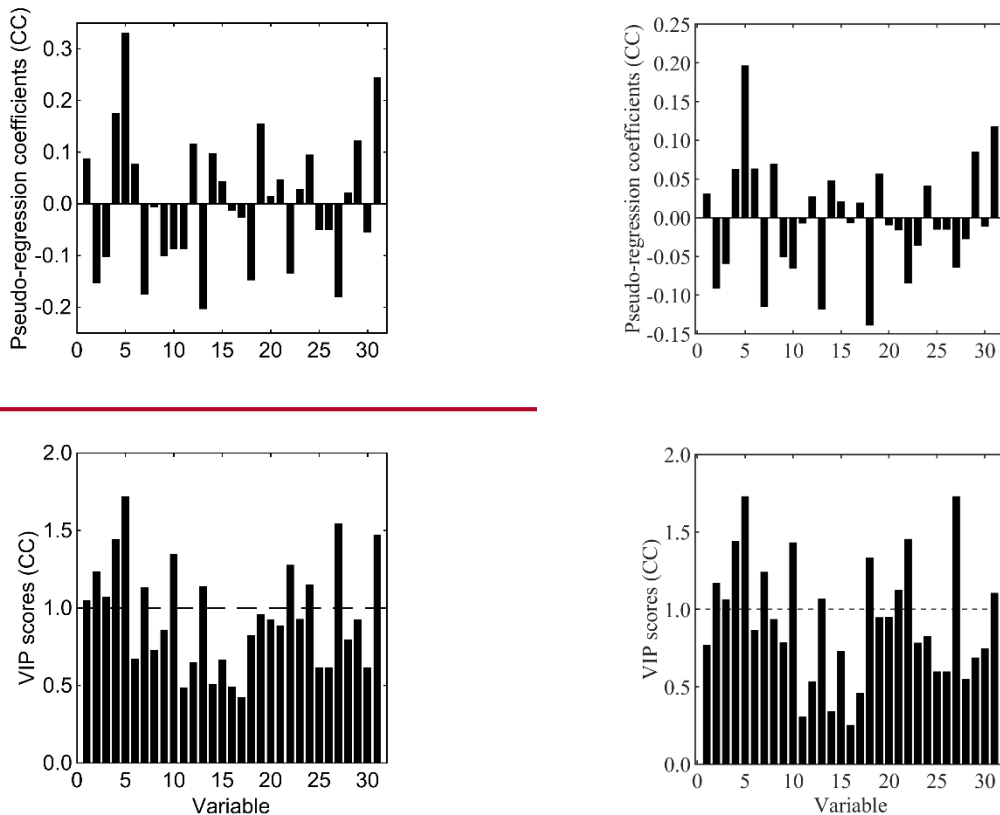
400 ~~The adipic acid/hexamethylenetetramine (11) CC sample was characterized by high leverage and high Hotelling's  $T^2$  value,~~  
401 ~~due to different behavior in terms of rotatable bonds and number of rings present in the structure. A similar behaviour, in terms~~  
402 ~~of difference in the number of rings, was observed also by two additional non-influential BM samples based on~~  
403 ~~hexamethylenetetramine paired with limonene (96) and menthone (103). In addition, these samples held high squared~~  
404 ~~residuals  $Q$  and, therefore, were characterized by features that did not align with the ones of the other samples.~~

405 ~~Finally, the ferulic acid/pyrazine (22) CC sample showed Hotelling's  $T^2$  values outside the 95% confidence interval due to~~  
406 ~~the different behaviour of the partner molecules in terms of molecular weight, connectivity, surface area, and electronic~~  
407 ~~properties, i.e., isotropic and anisotropic polarizability and number of valence electrons.~~

408 ~~Suspicious samples 22 (ferulic acid/pyrazine, CC) and 103 (menthone/hexamethylenetetramine, BM) appeared also to have~~  
409 ~~high residuals in absolute value and were wrongly assigned to their class in cross validation. Along with them, also three not~~  
410 ~~anomalous CC samples and one not anomalous BM samples, namely cinnamaldehyde/4-hydroxybenzoic acid (17),~~

411 carvacrol/nicotinamide (60), thymol/tetramethylpyrazine (68) and eugenol/pyrazine (92), respectively, were misclassified. In  
412 fact, their features were inversely related to their respective class. Nevertheless, the exclusion of the samples discussed above  
413 from the calibration set would not have produced any difference in terms of rotation of the LV space due to their low leverage.  
414 The correlation between class membership, coded in the  $Y_{cal}$ , and the predictors contained in the  $X_{cal}$  space can be observed  
415 in the PLS weights plot (Figure S3). ~~The variables involved in the discrimination are those whose weights follow the~~  
416 ~~discriminant direction: BM samples reach positive values of LV 1 and LV 2 in the  $Y_{cal}$  loading space (not shown), whereas it~~  
417 ~~is the opposite for CC samples. Therefore, it can be inferred that significant differences in descriptors related to polarizability,~~  
418 ~~exposed surface, and volume, such as atom and bond count (2, 3), molecular volume (9), octanol-water partition coefficient~~  
419 ~~(10), heteroatom count (27), topological polar surface area (22) are likely to prevent the formation of a cocrystal.~~  
420 Information regarding the contribution of each predictor involved in the discrimination of the modelled classes can be inferred  
421 by inspecting the pseudo-regression coefficients and the VIP score plots, reported in Figure 5. The latter parameter denotes  
422 the relative importance of each predictor of the  $X_{cal}$  space in the PLS-DA model in explaining the class membership encoded  
423 in the  $Y_{cal}$  and may guide variable selection. Generally, a variable can be considered important with a VIP score > 1; by  
424 contrast, a VIP score significantly lower than 1 indicates that a given variable is a good candidate for exclusion. According to  
425 the negative sign of the pseudo-regression coefficients related to CC class, ~~Therefore, it can be stated~~ inferred that significant  
426 differences in descriptors related to polarizability and exposed surface, such as atom and bond count (2, 3), octanol-water  
427 partition coefficient (10), topological polar surface area (22) and heteroatom count (27) are likely to prevent the formation of  
428 a cocrystal. It should be noted that this consideration agrees with what has emerged earlier from unsupervised modelling, and  
429 it is largely in agreement with widely applied rules of thumb in crystal engineering.

430



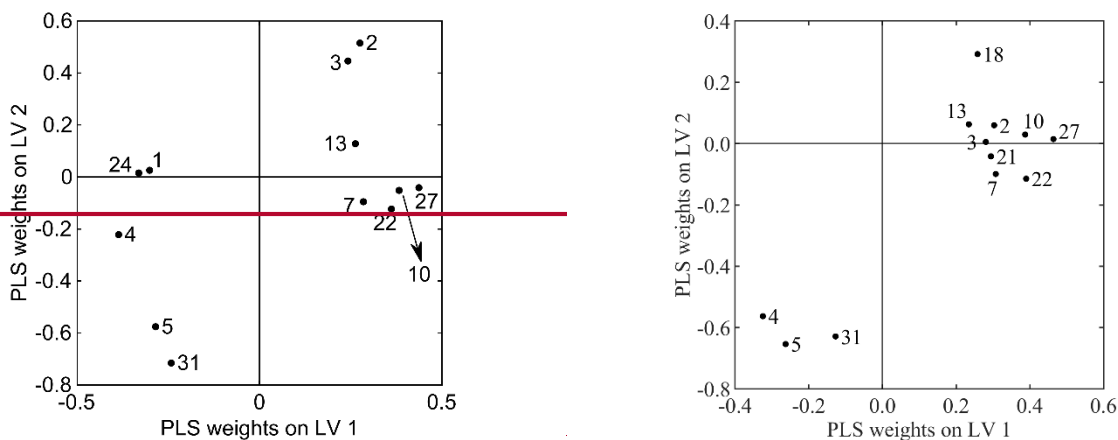
431

432 **Figure 5.** Pseudo-regression coefficients plot for the CC class (top). VIP scores related to each variable included in the  $X_{cal}$   
 433 space for the CC class, the significance threshold of 1 is depicted as a dashed horizontal line (bottom). Both plots report the  
 434 variable identification number on the  $x$  axis.

435

436 The variables that contributed most to the PLS weights were characterized by a VIP score > 1. Finally, a reduced PLS-DA  
 437 model based only on the important variables in agreement with the VIP scores was computed because of its ease of  
 438 interpretation. A 4 LVs model was calculated according to the maximum NERACC% in cross validation obtaining a  
 439 classification performance very similar to that achieved by including all descriptors. The weights plot of the reduced model  
 440 is provided in Figure 6, in which the weights of the variables involved in the discrimination are located on the positive and  
 441 negative sides of the first LV.

442



**Figure 6.** PLS weights of the variables for LV 1 vs. LV 2 of the reduced model.

Despite the interpretation of the correlation pattern among the variables being not so straightforward, some main considerations can be drawn. As a general comment, it can be stated that a good balance in the hydrogen bond propensity between the two partner molecules has to be achieved, being the difference in hydrogen bond acceptors (7) and number of heteroatoms (27) on positive PLS weights on LV 1. In addition, the coformers should have a similar behavior in terms of polarity, being the dissimilarities in polar surface area (22) and octanol/water partition coefficient (10) on positive weights on LV 1. It should be noted that this consideration agrees with what has emerged earlier from unsupervised modelling, and it is largely in agreement with widely applied rules of thumb in crystal engineering.

#### 4.2. Prediction of unknown samples Benchmarking on an external validation set

In order to evaluate the predictive capability of the model, an external set of  $N=58$  binary combinations of partner molecules was used. An overview of the involved samples is reported in Table S4 along with their estimated  $\hat{y}$  values.

This external validation set consists of 27 pairs classified experimentally as BMs and 31 pairs classified as CCs, i.e., 11 with our mechanochemistry/PXRD protocol and 20 retrieved from the CSD. The latter ones are reported in Section 3 of the Supplementary Information together with their CSD refcode and the reference to the original publications.

A graphical representation of the predicted values  $\hat{y}$  for the CC class is reported in Figure 7, as well as the squared residuals  $Q$  vs. Hotelling's  $T^2$  plot. Although there were some samples that did not conform to the model space, not all of them have been systematically misclassified. The confusion matrix is reported in Table 3 to summarize the results. In total, about 62.74% of the predictions were in agreement with the experimental results.

464

465 Table 3. Confusion matrix of the external validation set for the PLS-DA model. The second-last line shows the SEN% for the  
 466 modelled classes and last line shows the ACC%.

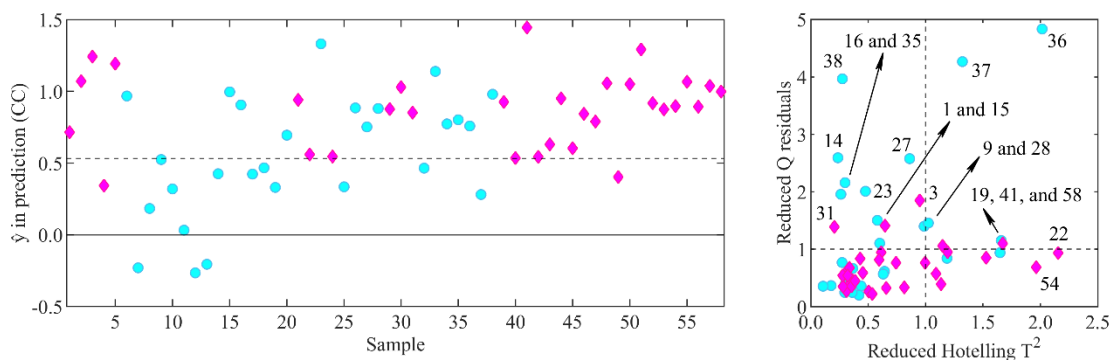
	<u>External validation set</u>	
	<u>Predicted as</u>	
	<u>CC</u>	<u>BM</u>
<u>Experimental CC</u>	<u>29</u>	<u>2</u>
<u>Experimental BM</u>	<u>13</u>	<u>14</u>
<u>SEN%</u>	<u>94%</u>	<u>52%</u>
<u>ACC%</u>	<u>74%</u>	

467

468 Specifically, 26-29 CC cases out of 31 were correctly classified. On the other hand, there were 17-13 false positive results and  
 469 10-14 cases in which the pairs were correctly identified as BMs. The model appears therefore to be quite conservative in  
 470 discarding the possibility of cocrystallization; hence, fewer potential new materials could be overlooked. ~~Similar behavior~~  
 471 ~~was observed also in the test set, meaning that the missed-discovery rate does not get worse when working on completely~~  
 472 ~~external data, thus demonstrating the stability of the model.~~ In addition, the fraction of misclassified CC samples (see Table  
 473 2 and Table 3) does not change significantly in the external validation set with respect to the test set.

474 As shown in Figure 7, Table 3 and Table S4, only 2 pairs of partner molecules predicted as BMs actually formed CCs, thus  
 475 producing 2 false negative results. It can be stated that these ~~Moreover, all the 5 false negative~~ results were borderline cases  
 476 when the estimated error [63] on the prediction was taken into consideration (data not shown): in fact their  $\hat{f}$  value was close  
 477 to the classification threshold of 0.53. The same is true for 8-6 of the false positive cases, while in the remaining 9-7 cases the  
 478 model confidently classified the pairs as CCs, in contrast to our experimental results. These findings could be ascribed to the  
 479 use of different preparation methods other than mechanochemical grinding. This hypothesis is somehow supported by the fact  
 480 that 19 out of the 20 CCs retrieved from the CSD, then prepared with a variety of methods, were correctly identified.

481



**Figure 7.**  $\hat{y}$  (CC samples) for the external validation set, the dashed horizontal line shows the hard classification threshold of 0.53 (left). Reduced squared residuals  $Q$  vs. reduced Hotelling's  $T^2$  plot, the dashed horizontal and vertical lines show the amplitude of the 95% confidence interval for both parameters (right). Samples are marked according to their class (magenta diamonds: cocrytals; cyan circles: binary mixtures).

## 5. Conclusion

This study highlighted the ability of a simple PLS-DA model to predict cocrytall formation without any *a priori* knowledge of the specific role of the involved partner molecules. Information deriving from both successful and unsuccessful cocrytallization experiments was used. The major advantage of the proposed methodology relies on the reduction of the experimental effort required for both the synthesis and characterization of new crystalline structures.

The model allows us to predict cocrytallization propensity with a ~~62~~74% of the predictions in agreement with the experimental results. Considering that the model was obtained on a training set spanning different molecular characteristics, it can be stated that it is suitable for a fairly general applicability.

Indeed, once in possess of the set of chemical descriptors for the molecules of interest, it is sufficient to calculate the absolute value of their difference and perform a linear combination using the pseudo-regression coefficients to obtain a prediction on cocrytall formation. The precalculated values for the set of descriptors comprising 2193 GRAS molecules are available in the Supplementary Material. By applying the proposed methodology, ~~seven~~ten new cocrytals were discovered and ~~an~~ additional ~~four~~ compounds ~~was~~were obtained by chance.



501 Another figure of merit of the proposed approach is the possibility of understanding through the inspection of PLS weights  
502 how the degree of similarity in terms of molecular features of the two partner molecules is correlated with the possibility of  
503 obtaining a cocrystal.

504 On a closing note, we would like to strongly encourage scientists to report failed attempts at cocrystallization along with the  
505 technique used, as access to this information could play a pivotal role in refining predictive models, making them less sensitive  
506 to selective reporting bias.

507

## 508 **Acknowledgements**

509 This work was funded by the Ministero delle Politiche Agricole, Alimentari, Forestali e del Turismo (MIPAAFT) by granting  
510 for the project PAC/Packaging Attivo Cristallino. This work has benefited from the framework of the COMP-HUB initiative,  
511 funded by the “Department of Excellence” program of the Italian Ministry for Education, University and Research (MIUR,  
512 2018–2022).

513

514

## 515 **Authorship Contribution Statement**

516 **Fabio Fornari:** Conceptualization (equal), Formal Analysis (lead), Methodology (equal), Validation (equal), Visualization  
517 (equal), Writing – Original Draft Preparation (lead), Writing – Review & Editing (equal).

518 **Fabio Montisci:** Investigation (equal), Data Curation (lead), Validation (equal), Visualization (equal), Writing – Original  
519 Draft Preparation (supporting), Writing – Review & Editing (equal).

520 **Federica Bianchi:** Conceptualization (equal), Methodology (lead), Resources (equal), Supervision (equal), Writing – Review  
521 & Editing (equal).

522 **Marina Cocchi:** Conceptualization (equal), Methodology (equal), Resources (equal), Software (lead), Writing – Review &  
523 Editing (equal).

524 **Claudia Carraro:** Investigation (equal), Data Curation (equal).

525 **Francesca Cavaliere:** Software (equal), Writing – Review & Editing (supporting).

526 **Pietro Cozzini:** Software (equal), Writing – Review & Editing (supporting).

527 **Francesca Peccati:** Software (supporting), Writing – Review & Editing (supporting).

528 **Paolo P. Mazzeo:** Conceptualization (equal), Investigation (supporting), Data Curation (supporting), Validation (supporting),  
529 Writing – Review & Editing (supporting).

530 **Nicolò Riboni:** Investigation (supporting), Writing – Review & Editing (supporting).

531 **Maria Careri:** Supervision (equal), Funding Acquisition (equal), Resources (equal), Writing – Review & Editing (equal).

532 **Alessia Bacchi:** Conceptualization (lead), Funding Acquisition (equal), Resources (equal), Supervision (equal), Writing –  
533 Review & Editing (equal).

534

535

### 536 **Declaration of Competing Interest**

537 The authors declare that they have no known competing financial interests or personal relationships that could have appeared  
538 to influence the work reported in this paper.

539 **References**

- 540 [1] E.V.R. Campos, P.L.F. Proença, J.L. Oliveira, M. Bakshi, P.C. Abhilash, L.F. Fraceto, Use of botanical insecticides  
541 for sustainable agriculture: Future perspectives, *Ecol. Indic.* 105 (2019) 483–495.  
542 doi:10.1016/J.ECOLIND.2018.04.038.
- 543 [2] Y. Kourkoutas, M. Angane, S. Swift, K. Huang, C.A. Butts, S.Y. Quek, Essential Oils and Their Major  
544 Components: An Updated Review on Antimicrobial Activities, Mechanism of Action and Their Potential  
545 Application in the Food Industry, *Foods* 2022, Vol. 11, Page 464. 11 (2022) 464. doi:10.3390/FOODS11030464.
- 546 [3] N.S. Singh, R. Sharma, T. Parween, P.K. Patanjali, Pesticide Contamination and Human Health Risk Factor, *Mod.*  
547 *Age Environ. Probl. Their Remediat.* (2018) 49–68. doi:10.1007/978-3-319-64501-8\_3.
- 548 [4] *Climate Change, Intercropping, Pest Control and Beneficial Microorganisms*, Springer Netherlands, 2009.  
549 doi:10.1007/978-90-481-2716-0.
- 550 [5] F. Bakkali, S. Averbeck, D. Averbeck, M. Idaomar, Biological effects of essential oils - A review, *Food Chem.*  
551 *Toxicol.* 46 (2008) 446–475. doi:10.1016/j.fct.2007.09.106.
- 552 [6] M. Alonso-Gato, G. Astray, J.C. Mejuto, J. Simal-Gandara, Essential Oils as Antimicrobials in Crop Protection,  
553 *Antibiotics.* 10 (2021) 34. doi:10.3390/antibiotics10010034.
- 554 [7] F. Bianchi, F. Fornari, N. Riboni, C. Spadini, C.S. Cabassi, M. Iannarelli, C. Carraro, P.P. Mazzeo, A. Bacchi, S.  
555 Orlandini, S. Furlanetto, M. Careri, Development of novel cocrystal-based active food packaging by a Quality by  
556 Design approach, *Food Chem.* 347 (2021) 129051. doi:10.1016/j.foodchem.2021.129051.
- 557 [8] A.T.H. Mossa, Green Pesticides: Essential oils as biopesticides in insect-pest management, *J. Environ. Sci. Technol.*  
558 9 (2016) 354–378. doi:10.3923/jest.2016.354.378.
- 559 [9] S. Sharma, S. Barkauskaite, A.K. Jaiswal, S. Jaiswal, Essential oils as additives in active food packaging, *Food*  
560 *Chem.* (2020) 128403. doi:10.1016/j.foodchem.2020.128403.
- 561 [10] Food and Drug Administration, 54960 Federal Register / Vol . 81 , No . 159 / Wednesday , August 17 , 2016 / Rules  
562 and Regulations, 81 (2016) 54960–55055.
- 563 [11] J. Wiczyńska, I. Cavoski, Antimicrobial, antioxidant and sensory features of eugenol, carvacrol and trans-anethole  
564 in active packaging for organic ready-to-eat iceberg lettuce, *Food Chem.* 259 (2018) 251–260.  
565 doi:10.1016/j.foodchem.2018.03.137.
- 566 [12] L. Pavoni, D.R. Perinelli, G. Bonacucina, M. Cespi, G.F. Palmieri, An overview of micro- and nanoemulsions as

- 567 vehicles for essential oils: Formulation, preparation and stability, *Nanomaterials*. 10 (2020) 135.  
568 doi:10.3390/nano10010135.
- 569 [13] G.R. Desiraju, *Crystal Engineering: A Holistic View*, *Angew. Chemie Int. Ed.* 46 (2007) 8342–8356.  
570 doi:10.1002/anie.200700534.
- 571 [14] A. Bacchi, P.P. Mazzeo, *Cocrystallization as a tool to stabilize liquid active ingredients*, *Crystallogr. Rev.* (2021) 1–  
572 22. doi:10.1080/0889311X.2021.1978079.
- 573 [15] D. Balestri, P.P. Mazzeo, R. Perrone, F. Fornari, F. Bianchi, M. Careri, A. Bacchi, P. Pelagatti, *Deciphering the*  
574 *Supramolecular Organization of Multiple Guests Inside a Microporous MOF to Understand their Release Profile*,  
575 *Angew. Chemie Int. Ed.* 60 (2021) 10194–10202. doi:10.1002/ANIE.202017105.
- 576 [16] D. Balestri, P.P. Mazzeo, C. Carraro, N. Demitri, P. Pelagatti, A. Bacchi, *Stepwise Evolution of Molecular*  
577 *Nanoaggregates Inside the Pores of a Highly Flexible Metal–Organic Framework*, *Angew. Chemie*. 131 (2019)  
578 17503–17511. doi:10.1002/ange.201907621.
- 579 [17] P.P. Mazzeo, S. Canossa, C. Carraro, P. Pelagatti, A. Bacchi, *Systematic coformer contribution to cocrystal*  
580 *stabilization: energy and packing trends*, *CrystEngComm*. 22 (2020) 7341–7349. doi:10.1039/D0CE00291G.
- 581 [18] P.P. Mazzeo, C. Carraro, A. Arns, P. Pelagatti, A. Bacchi, *Diversity through Similarity: A World of Polymorphs,*  
582 *Solid Solutions, and Cocrystals in a Vial of 4,4'-Diazopyridine*, *Cryst. Growth Des.* 20 (2020) 636–644.  
583 doi:10.1021/acs.cgd.9b01052.
- 584 [19] P.P. Mazzeo, C. Carraro, A. Monica, D. Capucci, P. Pelagatti, F. Bianchi, S. Agazzi, M. Careri, A. Raio, M. Carta,  
585 F. Menicucci, M. Belli, M. Michelozzi, A. Bacchi, *Designing a Palette of Cocrystals Based on Essential Oil*  
586 *Constituents for Agricultural Applications*, *ACS Sustain. Chem. Eng.* 7 (2019) 17929–17940.  
587 doi:10.1021/acssuschemeng.9b04576.
- 588 [20] D. Capucci, D. Balestri, P.P. Mazzeo, P. Pelagatti, K. Rubini, A. Bacchi, *Liquid Nicotine Tamed in Solid Forms by*  
589 *Cocrystallization*, *Cryst. Growth Des.* 17 (2017) 4958–4964. doi:10.1021/acs.cgd.7b00887.
- 590 [21] P.P. Mazzeo, M. Pioli, F. Montisci, A. Bacchi, P. Pelagatti, *Mechanochemical Preparation of Dipyridyl-*  
591 *Naphthalenediimide Cocrystals: Relative Role of Halogen-Bond and  $\pi$ - $\pi$  Interactions*, *Cryst. Growth Des.* 21  
592 (2021) 5687–5696. doi:10.1021/acs.cgd.1c00531.
- 593 [22] N.K. Duggirala, M.L. Perry, Ö. Almarsson, M.J. Zaworotko, *Pharmaceutical cocrystals: Along the path to improved*  
594 *medicines*, *Chem. Commun.* 52 (2016) 640–655. doi:10.1039/c5cc08216a.

- 595 [23] J.W. Steed, The role of co-crystals in pharmaceutical design., *Trends Pharmacol. Sci.* 34 (2013) 185–93.  
596 doi:10.1016/j.tips.2012.12.003.
- 597 [24] J.G.P. Wicker, L.M. Crowley, O. Robshaw, E.J. Little, S.P. Stokes, R.I. Cooper, S.E. Lawrence, Will they co-  
598 crystallize?, *CrystEngComm*. 19 (2017) 5336–5340. doi:10.1039/C7CE00587C.
- 599 [25] A. Bacchi, D. Capucci, M. Giannetto, M. Mattarozzi, P. Pelagatti, N. Rodriguez-Hornedo, K. Rubini, A. Sala,  
600 Turning Liquid Propofol into Solid (without Freezing It): Thermodynamic Characterization of Pharmaceutical  
601 Cocrystals Built with a Liquid Drug, *Cryst. Growth Des.* 16 (2016) 6547–6555. doi:10.1021/acs.cgd.6b01241.
- 602 [26] Y. Xiao, L. Zhou, H. Hao, Y. Bao, Q. Yin, C. Xie, Cocrystals of propylthiouracil and nutraceuticals toward  
603 sustained-release: Design, structure analysis, and solid-state characterization, *Cryst. Growth Des.* 21 (2021) 1202–  
604 1217. doi:10.1021/acs.cgd.0c01519.
- 605 [27] O. Shemchuk, S. D’Agostino, C. Fiore, V. Sambri, S. Zannoli, F. Grepioni, D. Braga, Natural Antimicrobials Meet  
606 a Synthetic Antibiotic: Carvacrol/Thymol and Ciprofloxacin Cocrystals as a Promising Solid-State Route to  
607 Activity Enhancement, *Cryst. Growth Des.* 20 (2020) 6796–6803. doi:10.1021/acs.cgd.0c00900.
- 608 [28] M.D. Perera, J. Desper, A.S. Sinha, C.B. Aakeröy, Impact and importance of electrostatic potential calculations for  
609 predicting structural patterns of hydrogen and halogen bonding, *CrystEngComm*. 18 (2016) 8631–8636.  
610 doi:10.1039/c6ce02089e.
- 611 [29] M.C. Etter, Encoding and Decoding Hydrogen-Bond Patterns of Organic Compounds, *Acc. Chem. Res.* 23 (1990)  
612 120–126. doi:10.1021/ar00172a005.
- 613 [30] C.A. Hunter, Quantifying intermolecular interactions: Guidelines for the molecular recognition toolbox, *Angew.*  
614 *Chemie - Int. Ed.* 43 (2004) 5310–5324. doi:10.1002/anie.200301739.
- 615 [31] M. Karimi-Jafari, L. Padrela, G.M. Walker, D.M. Croker, Creating cocrystals: A review of pharmaceutical cocrystal  
616 preparation routes and applications, *Cryst. Growth Des.* 18 (2018) 6370–6387. doi:10.1021/acs.cgd.8b00933.
- 617 [32] N. Issa, P.G. Karamertzanis, G.W.A. Welch, S.L. Price, Can the formation of pharmaceutical cocrystals be  
618 computationally predicted? I. Comparison of lattice energies, *Cryst. Growth Des.* 9 (2009) 442–453.  
619 doi:10.1021/cg800685z.
- 620 [33] M.A. Mohammad, A. Alhalaweh, S.P. Velaga, Hansen solubility parameter as a tool to predict cocrystal formation,  
621 *Int. J. Pharm.* 407 (2011) 63–71. doi:10.1016/j.ijpharm.2011.01.030.
- 622 [34] L. Fábíán, Cambridge Structural Database Analysis of Molecular Complementarity in Cocrystals, *Cryst. Growth*

623 Des. 9 (2009) 1436–1443. doi:10.1021/cg800861m.

624 [35] C.R. Groom, I.J. Bruno, M.P. Lightfoot, S.C. Ward, The Cambridge structural database, *Acta Crystallogr. Sect. B*  
625 *Struct. Sci. Cryst. Eng. Mater.* 72 (2016) 171–179. doi:10.1107/S2052520616003954.

626 [36] J. Devogelaer, H. Meekes, P. Tinnemans, E. Vlieg, R. Gelder, Co- crystal Prediction by Artificial Neural  
627 Networks\*\*, *Angew. Chemie Int. Ed.* 59 (2020) 21711–21718. doi:10.1002/anie.202009467.

628 [37] D. Wang, Z. Yang, B. Zhu, X. Mei, X. Luo, Machine-Learning-Guided Cocrystal Prediction Based on Large Data  
629 Base, *Cryst. Growth Des.* 20 (2020) 6610–6621. doi:10.1021/acs.cgd.0c00767.

630 [38] M. Przybyłek, T. Jeliński, J. Słabuszewska, D. Ziółkowska, K. Mroczyńska, P. Cysewski, Application of  
631 Multivariate Adaptive Regression Splines (MARSplines) Methodology for Screening of Dicarboxylic Acid  
632 Cocrystal Using 1D and 2D Molecular Descriptors, *Cryst. Growth Des.* 19 (2019) 3876–3887.  
633 doi:10.1021/acs.cgd.9b00318.

634 [39] M. Przybyłek, P. Cysewski, Distinguishing Cocrystals from Simple Eutectic Mixtures: Phenolic Acids as Potential  
635 Pharmaceutical Coformers, *Cryst. Growth Des.* 18 (2018) 3524–3534. doi:10.1021/acs.cgd.8b00335.

636 [40] A. Vriza, A.B. Canaj, R. Vismara, L.J. Kershaw Cook, T.D. Manning, M.W. Gaultois, P.A. Wood, V. Kurlin, N.  
637 Berry, M.S. Dyer, M.J. Rosseinsky, One class classification as a practical approach for accelerating  $\pi$ - $\pi$  co-crystal  
638 discovery, *Chem. Sci.* 12 (2021) 1702–1719. doi:10.1039/d0sc04263c.

639 [41] M.E. Mswahili, M.J. Lee, G.L. Martin, J. Kim, P. Kim, G.J. Choi, Y.S. Jeong, Cocrystal prediction using machine  
640 learning models and descriptors, *Appl. Sci.* 11 (2021) 1–12. doi:10.3390/app11031323.

641 [42] H. Moriwaki, Y.S. Tian, N. Kawashita, T. Takagi, Mordred: A molecular descriptor calculator, *J. Cheminform.* 10  
642 (2018) 4. doi:10.1186/s13321-018-0258-y.

643 [43] J.G.P. Wicker, R.I. Cooper, Will it crystallise? Predicting crystallinity of molecular materials, *CrystEngComm.* 17  
644 (2015) 1927–1934. doi:10.1039/c4ce01912a.

645 [44] I.J. Bruno, J.C. Cole, P.R. Edgington, M. Kessler, C.F. Macrae, P. McCabe, J. Pearson, R. Taylor, New software for  
646 searching the Cambridge Structural Database and visualizing crystal structures, *Acta Crystallogr. Sect. B Struct.*  
647 *Sci.* 58 (2002) 389–397. doi:10.1107/S0108768102003324.

648 [45] C.F. Macrae, I.J. Bruno, J.A. Chisholm, P.R. Edgington, P. McCabe, E. Pidcock, L. Rodriguez-Monge, R. Taylor, J.  
649 van de Streek, P.A. Wood, Mercury CSD 2.0 – new features for the visualization and investigation of crystal  
650 structures, *J. Appl. Crystallogr.* 41 (2008) 466–470. doi:10.1107/S0021889807067908.

- 651 [46] M. Baroni, G. Cruciani, S. Sciabola, F. Perruccio, J.S. Mason, A Common Reference Framework for  
652 Analyzing/Comparing Proteins and Ligands. Fingerprints for Ligands And Proteins (FLAP): Theory and  
653 Application, *J. Chem. Inf. Model.* 47 (2007) 279–294. doi:10.1021/ci600253e.
- 654 [47] Sybyl 8.1, Tripos International, St. Louis, USA, 2009.
- 655 [48] M. Clark, R.D. Cramer, N. Van Opdenbosch, Validation of the general purpose tripos 5.2 force field, *J. Comput.*  
656 *Chem.* 10 (1989) 982–1012. doi:10.1002/jcc.540100804.
- 657 [49] Schrödinger LLC, The PyMOL Molecular Graphics System, PyMOL Mol. Graph. Syst. Version 2.0. (2010).
- 658 [50] G. Landrum, RDKit: Open-Source Cheminformatics Software, [Http://Www.Rdkit.Org/](http://www.rdkit.org/). (2021).
- 659 [51] G.A. Frisch, M. J., Trucks, G. W., Schlegel, H. B., Scuseria, G. E., Robb, M. A., Cheeseman, J. R., Scalmani, G.,  
660 Barone, V., Petersson, Gaussian 16, Rev. A.03, Gaussian, Inc., Wallingford, CT. (2016).
- 661 [52] R.D. Snee, Validation of Regression Models: Methods and Examples, *Technometrics.* 19 (1977) 415–428.  
662 doi:10.1080/00401706.1977.10489581.
- 663 [53] M. Li Vigni, C. Durante, M. Cocchi, Exploratory Data Analysis, in: *Data Handl. Sci. Technol.*, Elsevier Ltd, 2013:  
664 pp. 55–126. doi:10.1016/B978-0-444-59528-7.00003-X.
- 665 [54] S. Wold, K. Esbensen, P. Geladi, Principal component analysis, *Chemom. Intell. Lab. Syst.* 2 (1987) 37–52.  
666 doi:10.1016/0169-7439(87)80084-9.
- 667 [55] R. Bro, A.K. Smilde, Principal component analysis, *Anal. Methods.* 6 (2014) 2812–2831. doi:10.1039/c3ay41907j.
- 668 [56] M. Cocchi, A. Biancolillo, F. Marini, Chemometric Methods for Classification and Feature Selection, in: *Compr.*  
669 *Anal. Chem.*, Elsevier B.V., 2018: pp. 265–299. doi:10.1016/bs.coac.2018.08.006.
- 670 [57] U.G. Indahl, H. Martens, T. Næs, From dummy regression to prior probabilities in PLS-DA, *J. Chemom.* 21 (2007)  
671 529–536. doi:10.1002/cem.1061.
- 672 [58] S. Wold, M. Sjöström, L. Eriksson, PLS-regression: A basic tool of chemometrics, in: *Chemom. Intell. Lab. Syst.*,  
673 Elsevier, 2001: pp. 109–130. doi:10.1016/S0169-7439(01)00155-1.
- 674 [59] D. Braga, E. Dichiarante, F. Grepioni, G.I. Lampronti, L. Maini, P.P. Mazzeo, S. D’Agostino, Mechanical  
675 Preparation of Crystalline Materials. An Oxymoron?, in: *Supramol. Chem.*, John Wiley & Sons, Ltd, Chichester,  
676 UK, 2012. doi:10.1002/9780470661345.smc115.
- 677 [60] J. Fürnkranz, P.K. Chan, S. Craw, C. Sammut, W. Uther, A. Ratnaparkhi, X. Jin, J. Han, Y. Yang, K. Morik, M.  
678 Dorigo, M. Birattari, T. Stützle, P. Brazdil, R. Vilalta, C. Giraud-Carrier, C. Soares, J. Rissanen, R.A. Baxter, I.

679 Bruha, R.A. Baxter, G.I. Webb, L. Torgo, A. Banerjee, H. Shan, S. Ray, P. Tadepalli, Y. Shoham, R. Powers, Y.  
680 Shoham, R. Powers, G.I. Webb, S. Ray, S. Scott, H. Blockeel, L. De Raedt, Manhattan Distance, *Encycl. Mach.*  
681 *Learn.* (2011) 639–639. doi:10.1007/978-0-387-30164-8\_506.

682 [61] L.E. Juarez-Orozco, O. Martinez-Manzanera, S. V. Nesterov, S. Kajander, J. Knuuti, The machine learning horizon  
683 in cardiac hybrid imaging, *Eur. J. Hybrid Imaging.* 2 (2018) 1–15. doi:10.1186/S41824-018-0033-3/FIGURES/5.

684 [62] A. Singh, N. Thakur, A. Sharma, A review of supervised machine learning algorithms, in: *3rd Int. Conf. Comput.*  
685 *Sustain. Glob. Dev.*, Institute of Electrical and Electronics Engineers, New Delhi, India, 2016: pp. 1310–1315.  
686 <https://ieeexplore.ieee.org/abstract/document/7724478>.

687 [63] N.M. Faber, R. Bro, Standard error of prediction for multiway PLS: 1. Background and a simulation study,  
688 *Chemom. Intell. Lab. Syst.* 61 (2002) 133–149. doi:10.1016/S0169-7439(01)00204-0.

689



## Highlights

A QSPR model for the discovery of cocrystals made by EOs and other GRAS molecules

Training set based on failed and successful cocrystallization experiments

Correct classification rate of 85% on the test set

Broad applicability and reduced experimental effort

**Declaration of interests**

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

## **Author Statement**

**Fabio Fornari:** Conceptualization (equal), Formal Analysis (lead), Methodology (equal), Validation (equal), Visualization (equal), Writing – Original Draft Preparation (lead), Writing – Review & Editing (equal).

**Fabio Montisci:** Investigation (equal), Data Curation (lead), Validation (equal), Visualization (equal), Writing – Original Draft Preparation (supporting), Writing – Review & Editing (equal).

**Federica Bianchi:** Conceptualization (equal), Methodology (lead), Resources (equal), Supervision (equal), Writing – Review & Editing (equal).

**Marina Cocchi:** Conceptualization (equal), Methodology (equal), Resources (equal), Software (lead), Writing – Review & Editing (equal).

**Claudia Carraro:** Investigation (equal), Data Curation (equal).

**Francesca Cavaliere:** Software (equal), Writing – Review & Editing (supporting).

**Pietro Cozzini:** Software (equal), Writing – Review & Editing (supporting).

**Francesca Peccati:** Software (supporting), Writing – Review & Editing (supporting).

**Paolo P. Mazzeo:** Conceptualization (equal), Investigation (supporting), Data Curation (supporting), Validation (supporting), Writing – Review & Editing (supporting).

**Nicolò Riboni:** Investigation (supporting), Writing – Review & Editing (supporting).

**Maria Careri:** Supervision (equal), Funding Acquisition (equal), Resources (equal), Writing – Review & Editing (equal).

**Alessia Bacchi:** Conceptualization (lead), Funding Acquisition (equal), Resources (equal), Supervision (equal), Writing – Review & Editing (equal).




Click here to access/download  
**Supplementary Material**  
Supplementary Material.docx



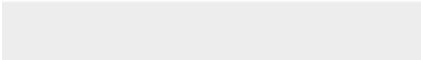



Click here to access/download  
**Supplementary Material**  
GRAS\_molecules\_dataset.xlsx





Click here to access/download  
**Supplementary Material**  
PLSDAmodel.m



1 **Chemometric-Assisted Cocrystallization: Supervised Pattern Recognition for Predicting the**  
2 **Formation of New Functional Cocrystals**

3

4 Fabio Fornari <sup>a</sup>, Fabio Montisci <sup>a</sup>, Federica Bianchi <sup>a,b,\*</sup>, Marina Cocchi <sup>c</sup>, Claudia Carraro <sup>a</sup>, Francesca Cavaliere <sup>d</sup>, Pietro  
5 Cozzini <sup>d</sup>, Francesca Peccati <sup>e</sup>, Paolo P. Mazzeo <sup>a,f</sup>, Nicolò Riboni <sup>a</sup>, Maria Careri <sup>a,g</sup>, Alessia Bacchi <sup>a,f</sup>

6

7 <sup>a</sup> *University of Parma, Department of Chemistry, Life Sciences and Environmental Sustainability, Parco Area delle Scienze*  
8 *17/A, 43124, Parma, Italy*

9 <sup>b</sup> *University of Parma, Interdepartmental Center for Packaging (CIPACK), Parco Area delle Scienze, 43124, Parma, Italy*

10 <sup>c</sup> *University of Modena and Reggio Emilia, Department of Chemical and Geological Sciences, Via Giuseppe Campi 103,*  
11 *41125, Modena, Italy*

12 <sup>d</sup> *University of Parma, Department of Food and Drug, Parco Area delle Scienze 17/A, 43124, Parma, Italy*

13 <sup>e</sup> *Basque Research and Technology Alliance (BRTA), Center for Cooperative Research in Biosciences (CIC bioGUNE),*  
14 *Bizkaia Technology Park 801A, 48160, Derio, Spain*

15 <sup>f</sup> *University of Parma, Biopharmanet-TEC, Parco Area delle Scienze 27/A, 43124, Parma, Italy*

16 <sup>g</sup> *University of Parma, Interdepartmental Center on Safety, Technologies, and Agri-Food Innovation (SITEIA.PARMA), Parco*  
17 *Area delle Scienze, 43124, Parma, Italy*

18

19

20

21

22

23

24 \* Author to whom correspondence should be addressed:

25 e-mail address: federica.bianchi@unipr.it; University of Parma, Department of Chemistry, Life Sciences and Environmental  
26 Sustainability, Parco Area delle Scienze 17/A, 43124, Parma, Italy

27 **Abstract**

28 Owing to the antimicrobial and insecticide properties, the use of natural compounds like essential oils and their active  
29 components has proven to be an effective alternative to synthetic chemicals in different fields ranging from drug delivery to  
30 agriculture and from nutrition to food preservation. Their limited application due to the high volatility and scarce water  
31 solubility can be expanded by using crystal engineering approaches to tune some properties of the active molecule by  
32 combining it with a suitable partner molecule (coformer). However, the selection of coformers and the experimental effort  
33 required for discovering cocrystals are the bottleneck of cocrystal engineering. This study explores the use of chemometrics  
34 to aid the discovery of cocrystals of active ingredients suitable for various applications. Partial Least Squares–Discriminant  
35 Analysis is used to discern cocrystals from binary mixtures based on the molecular features of the coformers. For the first  
36 time, by including failed cocrystallization data and considering a variety of chemically diverse compounds, the proposed  
37 method resulted in a successful prediction rate of 85% for the test set in the model validation phase and of 74% for the external  
38 test set.

39

40

41

42

43

44

45

46

47

48

49

50 **Keywords**

51 cocrystal, crystal engineering, chemoinformatics, chemometrics, partial least square discriminant analysis, Quantitative  
52 Structure–Property Relationship



53 **List of abbreviations (sorted alphabetically)**

54 ACC%; Accuracy

55 BM: Binary Mixture

56 CC: Cocrystal

57 CCDC: Cambridge Crystallographic Data Centre

58 CSD: Cambridge Structural Database

59 EO: Essential Oil

60 FDA: Food and Drug Administration

61 GRAS: Generally Recognized As Safe

62 LV: Latent Variable

63 MEP: Molecular Electrostatic Potential

64 PC: Principal Component

65 PCA: Principal Component Analysis

66 PLS-DA: Partial Least Squares-Discriminant Analysis

67 PXRD: Powder X-Ray Diffraction

68 QSPR: Quantitative Structure-Property Relationship

69 SEN%: Sensitivity

70 VIP: Variable Importance in Projection.

71

72

73

74

75

76

77

78

79

80

## 81 **1. Introduction**

82 In the last few decades, the use of agrochemicals and food preservatives has grown exponentially as a direct consequence of  
83 the rapid increase of the world population [1,2]. Owing to their potential adverse effects on both human health and  
84 environment [2–4], alternative strategies based on the use of more sustainable chemicals have been proposed to support the  
85 food system. Being able to exert antimicrobial, insecticidal, and antioxidant properties [5,6], essential oils (EOs) and their  
86 active components have been used as green substitutes of synthetic chemicals to extend the shelf-life of foodstuff and in pests  
87 control [7–9]. These compounds are Generally Recognized As Safe (GRAS) by the Food and Drug Administration (FDA)  
88 [10], however, despite their appealing properties, their use is limited by their high volatility and poor stability [7,9,11,12].

89 In fact, physicochemical properties of materials play a key role in determining whether a chemical is suitable for a specific  
90 purpose, thus strongly affecting its field of application. Scientists have always desired to obtain materials with target  
91 properties, and crystal engineering is one of the most interesting approaches to synthesize a great variety of crystalline  
92 materials for applications in various fields, ranging from pharmaceuticals to agrochemicals, and from nutraceuticals to  
93 cosmetics [13–16]. The basic idea of crystal engineering is related to the possibility of controlling the crystal structure of  
94 molecules and, therefore, the properties of the resulting solids. Polymorphism, vitrification and cocrystallization are some of  
95 the available strategies to modify the intrinsic properties of molecules without the need of synthetic modifications [17–21].

96 Cocrystals are multicomponent crystalline solid materials in which the constituents (i.e., coformers) are bound in a well-  
97 defined stoichiometric ratio [22,23] *via* non-covalent interactions (e.g., hydrogen bonds, halogen bonds,  $\pi$ - $\pi$  stacking) within  
98 the same crystal structure. Cocrystallization allows for the combination of the desired molecule of interest with properly  
99 selected partner molecules, paving the way to an array of potential materials with enhanced properties [19,24,25]. Within this  
100 frame of reference, cocrystals based on the active components of EOs have been proposed as active ingredients for food  
101 packaging, agrochemical and pharmaceutical applications [7,19,26,27].

102 Despite the great advantages offered by cocrystallization, the proper degree of complementarity between the two partner  
103 molecules required to obtain crystalline materials with the desired properties is not easy to assess [28–31]. In this context, the  
104 selection of coformers and the great effort required for both the systematic experimental screening and careful characterization  
105 of the products derived from the combination of all the possible coformer pairs represent the major bottleneck of cocrystal  
106 engineering. Computational techniques represent a powerful tool to reduce the experimental effort required for the discovery  
107 of new cocrystals, enabling to evaluate beforehand whether a cocrystal can be obtained starting from pre-selected coformers.  
108 These *in silico* strategies can be based on the calculation of a variety of parameters useful to predict the formation of a

109 cocrystal, such as lattice energy [32], solubility [33], hydrogen bond propensity along with the quantitation of molecular  
110 interaction energy [29,30], and molecular complementarity [34].

111 Despite the massive efforts spent to develop a method to predict cocrystal formation, at present none of the proposed strategies  
112 has proven to be both totally reliable and easy to apply.

113 Chemometrics could play a pivotal role in cocrystal discovery: up to now only few Machine Learning methods have been  
114 proposed in predicting cocrystal formation, enabling the screening of new cocrystals once a supervised model is properly  
115 trained and validated. In the study proposed by Devogelaer et al., information of successful cocrystallization experiments was  
116 directly taken from the Cambridge Structural Database (CSD) [35] and Artificial Neural Networks (ANN) were used to predict  
117 the formation of new cocrystals [36]. Similarly, Wang et al. relied on a consensus method based on multiple Random Forest  
118 algorithms, in which the successful cocrystallization dataset was integrated with randomly generated failed cocrystallization  
119 data [37]. These approaches are reported in the literature as network-based methods. Additional studies were based on the use  
120 of successful and unsuccessful cocrystallization datasets obtained from experimentation, literature, and/or the CSD for  
121 screening specific classes of cofomers. Within this framework, Przybyłek et al. used Multivariate Adaptive Regression  
122 Splines to predict the formation of dicarboxylic and phenolic acid-based cocrystals [38,39], whereas Wicker et al. focused on  
123 variously substituted benzoic acids and benzamides using a Support Vector Machine algorithm [24]. Vriza et al. used an  
124 ensemble one-class classification method to aid the discovery of  $\pi$ - $\pi$  cocrystals, thus giving a great contribution in enriching  
125 one of the most under-represented classes of cocrystals in the CSD [40]. Most recently, Mswahili et al. developed a cocrystal  
126 screening method based on ANN by using both successful and unsuccessful experimental cocrystallization data retrieved from  
127 the literature and a plethora of molecular descriptors calculated using Mordred [41,42].

128 In the frame of a research activity dealing with the synthesis of new functional cocrystals based on the active constituents of  
129 EOs and other GRAS molecules to broaden their applicability in the industrial field [7,19], we propose a chemometric  
130 approach to aid the discovery of new cocrystalline materials.

131 For the first time, a training set based on the results of failed (binary mixtures, BM) and successful (cocrystal, CC)  
132 cocrystallization experiments was used for the computation of a Quantitative Structure–Property Relationship-like (QSPR)  
133 model based on Partial Least Squares–Discriminant Analysis (PLS–DA), after preliminary exploratory analysis by Principal  
134 Component Analysis (PCA). The PLS-DA approach, with respect to network-based methods offers the advantages of having  
135 only one parameter to optimize, i.e., the model dimensionality, and the direct interpretation of the importance of descriptors  
136 in classification, while highlighting their interplay (by inspection of weights and loadings plots).

137 The effectiveness of the study relies on the use of compounds belonging to different chemical classes and a reduced number  
138 of 1D, 2D, and 3D molecular descriptors of various nature (e.g., constitutional, geometric, physical, topological, and surface  
139 area-based descriptors) [24,38,39,43], enabling the high-throughput screening of novel cocrystalline materials and offering  
140 maximum flexibility and effectiveness at a minimum computational and experimental cost.

141

## 142 **2. Experimental Procedures**

### 143 2.1. Mechanochemical protocol and class assignation

144 All the molecules in the dataset were chosen among the list of GRAS molecules drawn up by the FDA [10]. Selected pair of  
145 molecules among the chosen ones were assigned either to the CC class or to the BM class. Pairs of molecules in the dataset  
146 for which a cocrystal structure was already described in literature were individuated in the Cambridge Structural Database  
147 (CSD) [35] with the Cambridge Crystallographic Data Centre (CCDC) software ConQuest [44] and visualized with Mercury  
148 [45]. They are reported in Section 3 of the Supplementary Material, together with their unique CSD refcode and the reference  
149 to the original publications.

150 Cocrystallization for all the pairs with no known structure in literature was instead attempted with the following  
151 mechanochemical protocol. All the reagents employed were commercially available and used as such in all the experiments.  
152 Equimolar amounts of each reagent were directly mixed in an agate mortar and subjected to manual grinding for 10-15  
153 minutes, without using any solvent. The resulting powder samples were collected in closed vials. Assignation to CC or BM  
154 classes was performed by comparing the Powder X-ray Diffraction (PXRD) pattern of the ground sample with those of the  
155 pure reagents. Possible occurrence of polymorphic transitions for the reagents was excluded by comparing the experimental  
156 PXRD data after milling with the calculated pattern of all the known crystalline forms of the reagents. The occurrence of new  
157 peaks, unexplained by the presence of unreacted reagents, was taken as indication that cocrystallization had occurred and the  
158 sample was assigned to the CC class. In case no additional peaks appeared in the PXRD pattern the sample was instead  
159 classified as a BM.

160

161

### 162 2.2. Powder X-ray diffraction

163 Typically, PXRD data were collected on a Rigaku Smartlab XE diffractometer in  $\theta$ - $\theta$  Bragg-Brentano geometry with Cu K $\alpha$   
164 radiation. The samples were placed on glass supports and exposed to radiation ( $1.5^\circ \leq 2\theta \leq 50^\circ$ ) at a scan rate of  $10^\circ/\text{min}$ . The  
165 diffracted beam was collected on a 2D Hypix 3000 solid state detector.  $5^\circ$  radiant soller were used as a compromise for high  
166 flux and moderate peak asymmetry at low angles. Beam stopper and anti-scatterer air component were used to mitigate the  
167 profile at low angle. In some rare cases, the data were collected on a Thermo Fisher Scientific ARL X'TRA diffractometer in  
168  $\theta$ - $\theta$  Bragg-Brentano geometry with Cu K $\alpha$  radiation ( $3^\circ \leq 2\theta \leq 30^\circ$  at a scan rate of  $5^\circ/\text{min}$ , or  $3^\circ \leq 2\theta \leq 40^\circ$  at a scan rate of  
169  $0.3^\circ/\text{min}$ ).

170

### 171 **3. Computational Methods**

#### 172 **3.1. Molecular descriptors calculation**

173 For each molecule 31 molecular descriptors were calculated (Table S1). A theoretical background for the less known  
174 descriptors is given in Section 4 of the Supplementary Material.

175 The molecular weight, the number of atoms, the number of bonds, the number of hydrogen bond donor sites and the number  
176 of hydrogen bond acceptor sites were calculated with FLAP software (Fingerprint for Ligand and Protein) [46] at pH 7.0,  
177 using the 3D structures of all molecules in SDF format as input (downloaded from the PubChem database). The number of  
178 rotatable bonds, the number of rings, the hydrophobicity (accounted as the number of hydrophobic centers), the logP  
179 (logarithm of octanol/water partition coefficient), the molecular volume, the total molecular dipole moment (based on point  
180 charge distribution in the molecule), and its components along the axes (using the principal axes of the molecular graph) were  
181 then calculated for the same structures using Sybyl 8.1 [47] (www.tripos.com) and taking in consideration the protonation  
182 state of molecules. The same software was also used to estimate the strain energy of the molecule without performing any  
183 geometry optimization. This energy term relies on an electrostatic calculation from atomic charges using the internal Tripos  
184 force field [48]. For the estimation of molecular volume and dipole moment, a specific SPL script was employed. The  
185 calculated volume is enclosed in a water-accessible surface computed at a repulsive interaction energy of 0.20 kcal/mol with  
186 a water probe. A custom Python script was used to automatically calculate the Solvent Accessible Surface Area (SASA) in  
187 PyMol 2.0 [49], with the dot density parameter set to 4. The number of heteroatoms, the number of valence electrons, and the  
188 indexes  ${}^0\chi$ ,  ${}^0\chi^n$ ,  ${}^0\chi^v$ ,  $\alpha_{HK}$ ,  ${}^1\kappa_a$ , LabuteASA, SMR\_VSA, PEOE\_VSA, and TPSA were obtained running a Python 3.7 code with  
189 the open-source cheminformatics toolkit RDKit Q4 2013 [50]. The average isotropic polarizability  $\alpha_{iso}$ , the polarizability

190 anisotropy  $\Delta\alpha$ , and the Molecular Electrostatic Potential (MEP) were calculated with Gaussian 16 [51] following the *in vacuo*  
191 Density-Functional Theory optimization of all the molecules, employing the hybrid functional B3LYP and the People double-  
192 z basis set 6-31+g(d,p).  
193 Postprocessing of the MEP to extract critical points at a given electron density isosurface was performed with a custom Python  
194 3.6.1 script on a three-dimensional map (cube format) with a sampling density of 6 points/Bohr along the three directions.  
195 The MEP was analyzed at an electron density isosurface of 0.002 a.u. with a tolerance of 0.001 a.u., meaning that only MEP  
196 values corresponding to regions of space with electron density in the 0.001–0.003 a.u. range were considered. A first set of  
197 critical points was identified comparing MEP values of each cube point with those of its 6 nearest neighbors. A point was  
198 considered a local minimum if the number of nearest neighbors with higher MEP was greater or equal to a given integer (4).  
199 Likewise, a point was considered a local maximum if the number of nearest neighbors with lower MEP was greater or equal  
200 to the same integer. This first step yielded a large number of candidate critical points encompassing a wide range of MEP  
201 values. Since our focus was on identifying the regions of the molecules likely to be involved in strong hydrogen bonds within  
202 the cocrystal, in a second step this first set of points was filtered based on the MEP values of the global minimum and  
203 maximum. This was done as follows: for each local minimum (maximum), the ratio between its MEP value and that of the  
204 global minimum (maximum) was computed, and the point was kept only if the ratio exceeded a given threshold (0.1). In this  
205 way, only points corresponding to shallow critical points were discarded. This step allowed to identify the MEP isosurface  
206 regions corresponding to hydrogen bond donors and acceptors. However, due to the rugged character of the MEP map,  
207 multiple critical points of the same type could appear in close proximity. To univocally map a given region of the isosurface  
208 to a MEP value, critical points close to each other (below a distance threshold of 1.0 Bohr) were merged iteratively, keeping  
209 only the lower MEP point for minima and higher MEP point for maxima. The algorithm then provided the final set of MEP  
210 critical points at the given electron density isosurface.

211

## 212 3.2. Data analysis

213 The entire data analysis was carried out in MATLAB R2019a environment (Mathworks, Natick, Massachusetts, USA) with  
214 the aid of the PLS\_Toolbox 8.7.1 (Eigenvector Research Inc., Washington, USA) chemometric package.

215

216

### 217 3.2.1. Data preprocessing

218 Each sample was described by  $m = 31$  variables (Table S1): the absolute value of the difference between the molecular  
219 descriptors of the two partner molecules was calculated, thus obtaining the predictor matrix  $\mathbf{X}$  ( $181 \times 31$ ). The class  
220 membership was binary encoded (1: belonging to the class; 0: otherwise) in a dummy matrix  $\mathbf{Y}$  ( $181 \times 2$ ) with each column  
221 representing one of the two modelled classes. The dataset was split in two subsets by using the Duplex algorithm [52]: 70%  
222 of the data were used as calibration set,  $\mathbf{X}_{\text{cal}}$  ( $127 \times 31$ ) and  $\mathbf{Y}_{\text{cal}}$  ( $127 \times 2$ ), whereas the remaining 30% were used as test set,  
223  $\mathbf{X}_{\text{test}}$  ( $54 \times 31$ ) and  $\mathbf{Y}_{\text{test}}$  ( $54 \times 2$ ). The calibration set and the test set are reported in Table S2 and Table S3, respectively.  
224 Before carrying out both exploratory multivariate data analysis and the computation of the supervised model, the calibration  
225 matrix  $\mathbf{X}_{\text{cal}}$  was preprocessed column-wise by performing mean centering and scaling to unit variance. Mean centering was  
226 applied on the response matrix  $\mathbf{Y}_{\text{cal}}$  to ensure the stability of the model.

227

### 228 3.2.2. *Exploratory multivariate data analysis*

229 PCA [53–55] was carried out preliminarily on the calibration set  $\mathbf{X}_{\text{cal}}$  to assess the distribution of the samples and to check for  
230 potential data structures. Reduction of data dimensionality is carried out through the linear combination of the original  
231 variables in a set of orthogonal ones, i.e., Principal Components (PCs), which identify the direction of maximum variance.  
232 This is summarized in the decomposition equation:

$$233 \mathbf{X}_{\text{cal}} = \mathbf{TP}^T + \mathbf{E}$$

234 where  $\mathbf{T}$  and  $\mathbf{P}$  represent, respectively, the coordinates of the samples projected in the reduced space, i.e., the scores, and the  
235 weights each original variable has on a given PC, i.e., the loadings. The deviations from the model are accounted in the error  
236 matrix  $\mathbf{E}$ .

237

### 238 3.2.3. *Supervised pattern recognition*

239 PLS–DA [56,57] was used to discriminate pairs of partner molecules whose combination forms CCs from the ones giving  
240 BMs. PLS–DA is based on PLS regression [58]. Briefly, this supervised technique decomposes the predictor matrix  $\mathbf{X}_{\text{cal}}$  and  
241 the dependent variables matrix  $\mathbf{Y}_{\text{cal}}$  in a PCA-like way and imposes inner linear relationships between the  $\mathbf{X}_{\text{cal}}$  and the  $\mathbf{Y}_{\text{cal}}$   
242 scores as follows:

$$243 \mathbf{U} = \mathbf{bT}$$

244 where  $T$  and  $U$ , are the  $X_{\text{cal}}$  and the  $Y_{\text{cal}}$  scores, respectively. This is accomplished by rotating the Latent Variable (LV) space  
 245 of  $X_{\text{cal}}$  through a weight matrix  $W$  in a way that maximizes the covariance between  $T$  and  $U$ . The PLS regression model is  
 246 summarized as:

$$247 \quad Y_{\text{cal}} = X_{\text{cal}}\mathbf{B} + \mathbf{E}$$

248 where  $\mathbf{E}$  is the error matrix and  $\mathbf{B}$  is the pseudo-regression coefficient matrix expressed according to the following equation:

$$249 \quad \mathbf{B} = \mathbf{W}(\mathbf{P}^T\mathbf{W})^{-1} \text{diag}(\mathbf{b}) \mathbf{Q}$$

250 where  $\mathbf{P}$  and  $\mathbf{Q}$  are the  $X_{\text{cal}}$  and the  $Y_{\text{cal}}$  loadings, respectively.

251 In this case, the dependent variables in the  $Y_{\text{cal}}$  matrix are defined as dummy variables, one for each modelled class, taking  
 252 values of 1 if the sample belongs to the class and 0 otherwise. Current implementation of PLS–DA may differ on the basis of  
 253 how the classification rule is defined. In this work, a pure discriminant rule (samples are assigned univocally to only one  
 254 category) was applied, and thus a sample is assigned to the class for which the predicted response  $\hat{y}$  is the highest (i.e.,  $\hat{Y}$   
 255 values are continuous and not dummy as they were codified).

256 The proper number of LVs was chosen according to the maximum accuracy (ACC%; i.e., the percentage of samples correctly  
 257 assigned to the respective class) in leave-more-out cross validation, adopting a Venetian blinds cancellation scheme with 10  
 258 splits (blind thickness: 1). This operation was carried out by running a custom MATLAB routine. The performance of the  
 259 classification model was evaluated both on the calibration and the test sets in terms of ACC% as well as showing the confusion  
 260 matrix. In addition, the sensitivity (SEN%, i.e., the percentage of samples within a class that were correctly assigned to their  
 261 class) was calculated for both classes.

262 The importance of each predictor was estimated in terms of Variable Importance in Projection (VIP score) [56]. The VIP  
 263 score of the  $j^{\text{th}}$  variable in the  $X$  space is defined as the component-wise sum of its PLS weight  $w_j$  on the  $f^{\text{th}}$  component  
 264 multiplied by the fraction of variance of the  $Y$  explained by that component, according to the following equation:

$$265 \quad VIP_j^2 = \frac{1}{SSY_{\text{tot}F}} \sum_{f=1}^F w_{jf}^2 SSY_f J$$

266 where  $J$  is the number of variables in the  $X$  space and  $F$  is the number of LVs that were retained. Since:

$$267 \quad \sum_{j=1}^J VIP_j^2 = J$$

268 the proposed threshold for determining whether a variable could be considered important is set to 1..



269 Finally, the predictive capability of the model was evaluated on an external set of  $N = 58$  binary combinations of partner  
270 molecules. An overview of the involved samples is reported in Table S4 along with their estimated  $\hat{y}$  values.

271 The number of entries in the BM and CC classes for training, test, and external validation sets is reported in Table 1.

272

273 **Table 1.** Number of CC and BM samples in the calibration, test, and external validation set. The last column reports the total  
274 number  $N$  of samples *per* set.

	CC	BM	$N$
Calibration set	71	56	127
Test set	30	24	54
External validation set	31	27	58

275

## 276 4. Results and Discussion

277 The cocrystallization experiments were carried out mechanochemically by manual neat grinding of the two substances. This  
278 method was selected among many possible others due to its simplicity and promptness, allowing us to screen several molecular  
279 pairs in a standardized way [21,59]. The classification as BM or CC was assessed by PXRD patterns (available in Section 3  
280 of the Supplementary Material). Cocrystals already present in the CSD were also included in our dataset.

281

### 282 4.1. Data analysis

283 Data handling prior to analysis can affect the way the model is trained, thus having consequences on its interpretation.

284 Since samples are the result of the combination of two partner molecules, each one described by its own set of descriptors,  
285 the concatenation strategy, i.e., listing the descriptors of the first coformer followed by those of the second coformer, was  
286 discarded due to the lack of commutation between the two sets of descriptors. In fact, the order in which the molecular  
287 descriptors are listed represents an *a priori* decision on which compound is acting as molecule of interest or partner molecule.

288 In our case this would be sub-optimal since many of the molecules in our dataset could assume both roles. Therefore, in order  
289 to address the problem described above a commutative strategy capable of avoiding the production of indeterminate forms  
290 should be chosen. Considering that in the dataset many constitutional molecular descriptors were characterized by a few non-  
291 zero values, the calculation of both products and ratios between the descriptors was discarded since additional zero values or  
292 indeterminate forms could be generated. Also, the division is non-commutative.

293 In order to combine the two partner molecules without imposition on their role, the absolute value of the difference between  
294 the molecular descriptors of the partner molecules was calculated and used to describe each binary combination [38,39],  
295 giving maximum flexibility to the model. The information achieved is still relevant and easily interpretable since it is related  
296 to the dissimilarity between the descriptors. In fact, absolute value of the differences between descriptors for each case are  
297 the elements of the Manhattan distance [60], one of the possible indexes to account for the dissimilarity between cases in a  
298 multivariate way.

299 In the present study, the use of basic linear modelling methods was preferred with respect to non-linear modelling, such as  
300 ANN [61,62], to keep the calculations as simple as possible, and to ensure a certain degree of interpretability of the results.  
301 Furthermore, the number of samples is too limited to ensure proper tuning of the ANN hyperparameters.

302 For a fruitful discussion, samples and variables are reported in the text according to their identification number as follows: i)  
303 samples are written in plain text; ii) variables are underlined>. The key is available in Tables S1–S3.

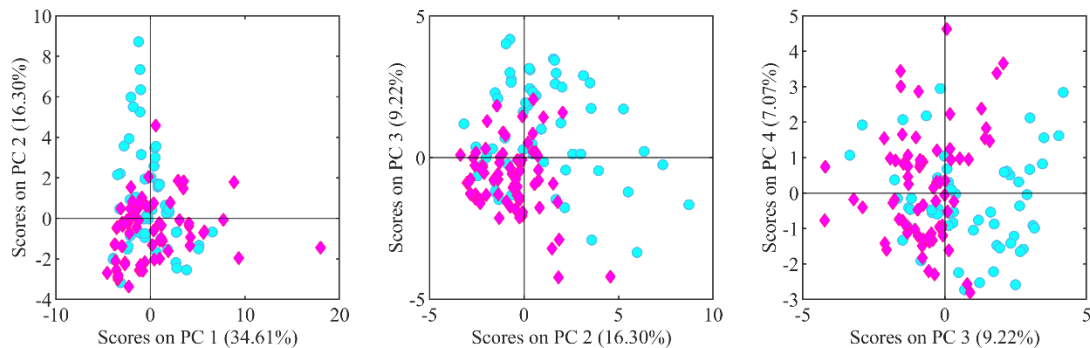
304

#### 305 4.1.1. Exploratory multivariate data analysis

306 After data preprocessing, PCA was used in an exploratory way to assess the presence of potential data structures in the  
307 calibration set.

308 Four PCs were retained explaining 67% of the variance. As shown in the score plots (Figure 1), a mild separation was present  
309 in the PC 2 vs. PC 3 score plot, with the groups separated by the bisecting line of the II and the IV quadrant. In addition, most  
310 of the CC samples occupied the III quadrant of the PC 2 vs. PC 3 score plot and were, in general, less scattered than BM  
311 samples.

312

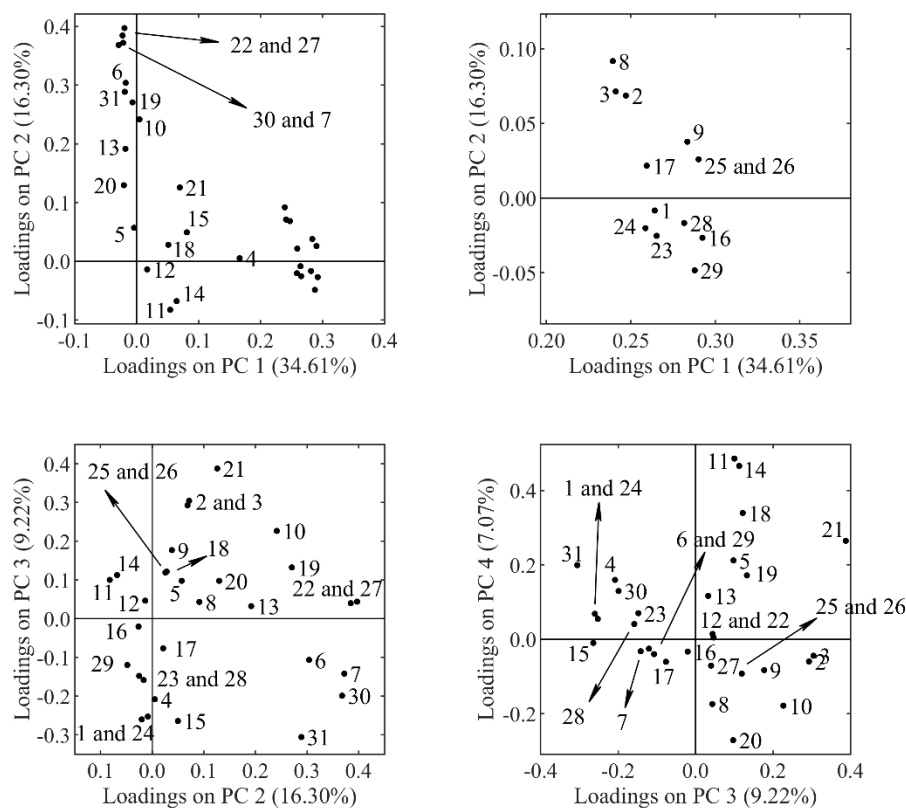


313

314 **Figure 1.** Scores of the samples on the first 4 PCs: PC 1 vs. PC 2 (left), PC 2 vs. PC 3 (center), PC 3 vs. PC 4 (right). The  
 315 fraction of variance explained by a given component is reported as a percentage value in parenthesis on the corresponding  
 316 axis. Samples are marked according to their class (magenta diamonds: cocrystals; cyan circles: binary mixtures).

317

318 As for the loading plots depicted in Figure 2, it can be observed that PC 1 explains the features related both to differences in  
 319 molecular dimensions (e.g., 1, 2, 3, 9, 16, 29) and connectivity (23, 24, 25, 26). PC 2 considers the dissimilarities in the  
 320 electronic properties (e.g., 6, 7, 19, 22, 30) of the two partner molecules as well as the difference in their number of  
 321 heteroatoms (27). PC 3 accounts for more specific features, such as differences in molecular complexity (21), molecular  
 322 refractivity (31), and energy (15). Lastly, PC 4 considers the dissimilarity in: i) component of the dipole along the  $x$  axis (11),  
 323 ii) total dipole (14) and iii) minimum of the MEP surface (20). Regarding the samples belonging to the CC class, these  
 324 compounds were characterized by partner molecules with a similar behavior in terms of molecular complexity (21) and  
 325 octanol/water partition coefficient (10).



326

327 **Figure 2.** Loading plots related to the PCA decomposition. PC 1 vs. PC 2 (top-left), magnification of the PC 1 vs. PC 2 (top-  
328 right), PC 2 vs. PC 3 (bottom-left), PC 3 vs. PC 4 (bottom-right). The fraction of variance explained by a given component is  
329 reported as a percentage value in parenthesis on the corresponding axis.

330

#### 331 4.1.2. Supervised pattern recognition

332 The relationship between the class membership and the variables was exploited by means of PLS–DA. Six LVs were retained  
333 according to maximum ACC% in cross validation. The PLS–DA model captured the 72% and 69% of the variance of the  $X_{cal}$   
334 and  $Y_{cal}$ , respectively. The values of the predicted response  $\hat{y}$  in cross validation related to CC samples are plotted in Figure  
335 S1.

336 A summary of the performance of the obtained model is reported in the confusion matrix (Table 2), whereas a graphical  
337 representation of the estimated and predicted values  $\hat{y}$  for the CC class is reported in Figure S2. As reported in Table 2, all the  
338 BM samples belonging to the test set were correctly classified except for 5 samples, whereas only 3 out of 30 CC samples  
339 were wrongly assigned to the BM class, obtaining a ACC% of 85%. Similarly, a high ACC% of 92% was obtained when the  
340 samples belonging to the calibration set were predicted by the model. The achieved results are extremely satisfactory, allowing  
341 for the *a priori* selection of the partner molecules required for the synthesis of novel cocrystals.

342

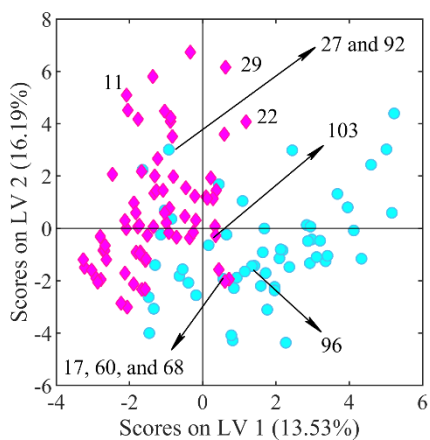
343 **Table 2.** Confusion matrix of the calibration and the test sets for the PLS–DA model. The second-last line shows the SEN%  
344 for the modelled classes and last line shows the ACC% for the calibration and the test sets.

	Calibration set		Test set	
	Predicted as		Predicted as	
	CC	BM	CC	BM
True CC	66	5	27	3
True BM	5	51	5	19
SEN%	93%	91%	90%	79%
ACC%	92%		85%	

345

346 The distribution of the samples in the reduced space of the LVs can be observed by inspecting the score plot (Figure 3),  
347 whereas information regarding suspicious and/or influential samples can be retrieved by the squared residuals  $Q$  vs.  
348 Hotelling's  $T^2$  and the residuals vs. leverage plots (Figure 4).

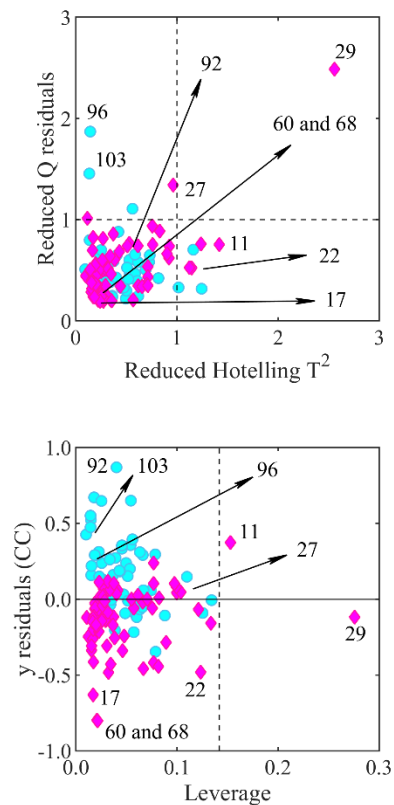
349



350

351 **Figure 3.** Distribution of the samples in the score space (LV 1 vs. LV 2). The fraction of variance explained by a given  
352 component is reported as a percentage value in parenthesis on the corresponding axis. Samples are marked according to their  
353 class (magenta diamonds: cocrystals; cyan circles: binary mixtures).

354



355

356 **Figure 4.** Reduced squared residuals  $Q$  vs. reduced Hotelling's  $T^2$  plot, the dashed horizontal and vertical lines show the  
 357 amplitude of the 95% confidence interval for both parameters (top). Residuals of CC samples vs. leverage plot, the dashed  
 358 vertical line shows the leverage limit (bottom). Samples are marked according to their class (magenta diamonds: cocrystals;  
 359 cyan circles: binary mixtures).

360

361 The maximum separation in the score space was provided by the first two LVs, with the CC samples located mainly at negative  
 362 scores on LV 1. By contrast, BM samples were more scattered and localized mostly on positive scores on LV 1.

363 A peculiar behavior was observed for samples 27 and 29. Sample 27 showed squared residual  $Q$  outside the 95% confidence  
 364 interval, whereas sample 29 showed both high Hotelling's  $T^2$  and high squared residuals  $Q$  together with high leverage in the  
 365  $Y_{cal}$  space. These CC samples were obtained by pairing fatty acids (lauric acid, 27, and palmitic acid, 29) with a low-molecular  
 366 weight coformer, i.e., nicotinamide. The variables responsible for this behavior can be related to the discrepancy in molecular  
 367 dimensions between the two partner molecules.

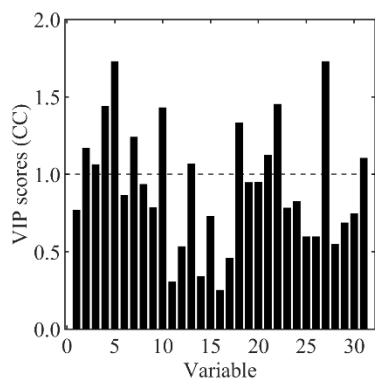
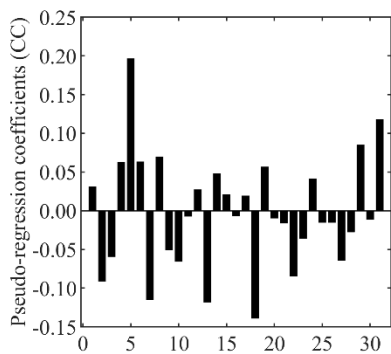
368 The adipic acid/hexamethylenetetramine (11) CC sample was characterized by high leverage and high Hotelling's  $T^2$  value,  
369 due to different behavior in terms of rotatable bonds and number of rings present in the structure. A similar behaviour, in terms  
370 of difference in the number of rings, was observed also by two additional non-influential BM samples based on  
371 hexamethylenetetramine paired with limonene (96) and menthone (103). In addition, these samples held high squared  
372 residuals  $Q$  and, therefore, were characterized by features that did not align with the ones of the other samples.

373 Finally, the ferulic acid/pyrazine (22) CC sample showed Hotelling's  $T^2$  values outside the 95% confidence interval due to  
374 the different behaviour of the partner molecules in terms of molecular weight, connectivity, surface area, and electronic  
375 properties, i.e., isotropic and anisotropic polarizability and number of valence electrons.

376 Suspicious samples 22 (ferulic acid/pyrazine, CC) and 103 (menthone/hexamethylenetetramine, BM) appeared also to have  
377 high residuals in absolute value and were wrongly assigned to their class in cross validation. Along with them, also three not  
378 anomalous CC samples and one not anomalous BM samples, namely cinnamaldehyde/4-hydroxybenzoic acid (17),  
379 carvacrol/nicotinamide (60), thymol/tetramethylpyrazine (68) and eugenol/pyrazine (92), respectively, were misclassified. In  
380 fact, their features were inversely related to their respective class. Nevertheless, the exclusion of the samples discussed above  
381 from the calibration set would not have produced any difference in terms of rotation of the LV space due to their low leverage.

382 The correlation between class membership, coded in the  $Y_{\text{cal}}$ , and the predictors contained in the  $X_{\text{cal}}$  space can be observed  
383 in the PLS weights plot (Figure S3).

384 Information regarding the contribution of each predictor involved in the discrimination of the modelled classes can be inferred  
385 by inspecting the pseudo-regression coefficients and the VIP score plots, reported in Figure 5. The latter parameter denotes  
386 the relative importance of each predictor of the  $X_{\text{cal}}$  space in the PLS-DA model in explaining the class membership encoded  
387 in the  $Y_{\text{cal}}$  and may guide variable selection. Generally, a variable can be considered important with a VIP score  $> 1$ ; by  
388 contrast, a VIP score significantly lower than 1 indicates that a given variable is a good candidate for exclusion. According to  
389 the negative sign of the pseudo-regression coefficients related to CC class, it can be stated that significant differences in  
390 descriptors related to polarizability and exposed surface, such as atom and bond count (2, 3), octanol-water partition  
391 coefficient (10), topological polar surface area (22) and heteroatom count (27) are likely to prevent the formation of a  
392 cocrystal. It should be noted that this consideration agrees with what has emerged earlier from unsupervised modelling, and  
393 it is largely in agreement with widely applied rules of thumb in crystal engineering.



394

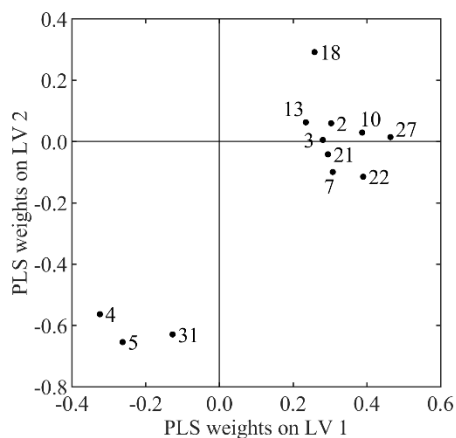
395 **Figure 5.** Pseudo-regression coefficients plot for the CC class (top). VIP scores related to each variable included in the  $X_{cal}$   
 396 space for the CC class, the significance threshold of 1 is depicted as a dashed horizontal line (bottom). Both plots report the  
 397 variable identification number on the  $x$  axis.

398

399 Finally, a reduced PLS-DA model based only on the important variables in agreement with the VIP scores was computed  
 400 because of its ease of interpretation. A 4 LVs model was calculated according to the maximum ACC% in cross validation  
 401 obtaining a classification performance very similar to that achieved by including all descriptors. The weights plot of the  
 402 reduced model is provided in Figure 6, in which the weights of the variables involved in the discrimination are located on the  
 403 positive and negative sides of the first LV.

404





405

406 **Figure 6.** PLS weights of the variables for LV 1 vs. LV 2 of the reduced model.

407

#### 408 4.2. Benchmarking on an external validation set

409 The external validation set consists of 27 pairs classified experimentally as BMs and 31 pairs classified as CCs, i.e., 11 with  
 410 our mechanochemistry/PXRD protocol and 20 retrieved from the CSD. The latter ones are reported in Section 3 of the  
 411 Supplementary Information together with their CSD refcode and the reference to the original publications.

412 A graphical representation of the predicted values  $\hat{y}$  for the CC class is reported in Figure 7, as well as the squared residuals  
 413  $Q$  vs. Hotelling's  $T^2$  plot. Although there were some samples that did not conform to the model space, not all of them have  
 414 been systematically misclassified. The confusion matrix is reported in Table 3 to summarize the results. In total, about 74%  
 415 of the predictions were in agreement with the experimental results.

416

417 **Table 3.** Confusion matrix of the external validation set for the PLS-DA model. The second-last line shows the SEN% for the  
 418 modelled classes and last line shows the ACC%.

	External validation set	
	Predicted as	
	CC	BM
Experimental CC	29	2
Experimental BM	13	14
SEN%	94%	52%

---

ACC%	74%
------	-----

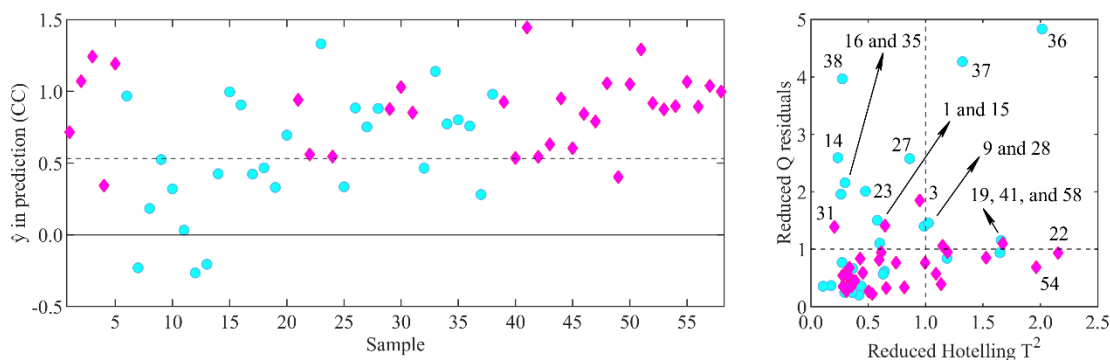
---

419

420 Specifically, 29 CC cases out of 31 were correctly classified. On the other hand, there were 13 false positive results and 14  
421 cases in which the pairs were correctly identified as BMs. The model appears therefore to be quite conservative in discarding  
422 the possibility of cocrystallization; hence, fewer potential new materials could be overlooked. In addition, the fraction of  
423 misclassified CC samples (see Table 2 and Table 3) does not change significantly in the external validation set with respect  
424 to the test set.

425 As shown in Figure 7, Table 3 and Table S4, only 2 pairs of partner molecules predicted as BMs actually formed CCs, thus  
426 producing 2 false negative results. It can be stated that these results were borderline cases when the estimated error [63] on  
427 the prediction was taken into consideration (data not shown): in fact their  $\hat{y}$  value was close to the classification threshold of  
428 0.53. The same is true for 6 of the false positive cases, while in the remaining 7 cases the model confidently classified the  
429 pairs as CCs, in contrast to our experimental results. These findings could be ascribed to the use of different preparation  
430 methods other than mechanochemical grinding. This hypothesis is somehow supported by the fact that 19 out of the 20 CCs  
431 retrieved from the CSD, then prepared with a variety of methods, were correctly identified.

432



433

434 **Figure 7.**  $\hat{y}$  (CC samples) for the external validation set, the dashed horizontal line shows the hard classification threshold of  
435 0.53 (left). Reduced squared residuals  $Q$  vs. reduced Hotelling's  $T^2$  plot, the dashed horizontal and vertical lines show the  
436 amplitude of the 95% confidence interval for both parameters (right). Samples are marked according to their class (magenta  
437 diamonds: cocrystals; cyan circles: binary mixtures).

438

## 439 **5. Conclusion**

440 This study highlighted the ability of a simple PLS–DA model to predict cocrystal formation without any *a priori* knowledge  
441 of the specific role of the involved partner molecules. Information deriving from both successful and unsuccessful  
442 cocrystallization experiments was used. The major advantage of the proposed methodology relies on the reduction of the  
443 experimental effort required for both the synthesis and characterization of new crystalline structures.

444 The model allows us to predict cocrystallization propensity with a 74% of the predictions in agreement with the experimental  
445 results. Considering that the model was obtained on a training set spanning different molecular characteristics, it can be stated  
446 that it is suitable for a fairly general applicability.

447 Indeed, once in possess of the set of chemical descriptors for the molecules of interest, it is sufficient to calculate the absolute  
448 value of their difference and perform a linear combination using the pseudo-regression coefficients to obtain a prediction on  
449 cocrystal formation. The precalculated values for the set of descriptors comprising 2193 GRAS molecules are available in the  
450 Supplementary Material. By applying the proposed methodology, ten new cocrystals were discovered and an additional  
451 compound was obtained by chance.

452 Another figure of merit of the proposed approach is the possibility of understanding through the inspection of PLS weights  
453 how the degree of similarity in terms of molecular features of the two partner molecules is correlated with the possibility of  
454 obtaining a cocrystal.

455 On a closing note, we would like to strongly encourage scientists to report failed attempts at cocrystallization along with the  
456 technique used, as access to this information could play a pivotal role in refining predictive models, making them less sensitive  
457 to selective reporting bias.

458

## 459 **Acknowledgements**

460 This work was funded by the Ministero delle Politiche Agricole, Alimentari, Forestali e del Turismo (MIPAAFT) by granting  
461 for the project PAC/Packaging Attivo Cristallino. This work has benefited from the framework of the COMP-HUB initiative,  
462 funded by the “Department of Excellence” program of the Italian Ministry for Education, University and Research (MIUR,  
463 2018–2022).

464

465

466 **Authorship Contribution Statement**

467 **Fabio Fornari:** Conceptualization (equal), Formal Analysis (lead), Methodology (equal), Validation (equal), Visualization  
468 (equal), Writing – Original Draft Preparation (lead), Writing – Review & Editing (equal).

469 **Fabio Montisci:** Investigation (equal), Data Curation (lead), Validation (equal), Visualization (equal), Writing – Original  
470 Draft Preparation (supporting), Writing – Review & Editing (equal).

471 **Federica Bianchi:** Conceptualization (equal), Methodology (lead), Resources (equal), Supervision (equal), Writing – Review  
472 & Editing (equal).

473 **Marina Cocchi:** Conceptualization (equal), Methodology (equal), Resources (equal), Software (lead), Writing – Review &  
474 Editing (equal).

475 **Claudia Carraro:** Investigation (equal), Data Curation (equal).

476 **Francesca Cavaliere:** Software (equal), Writing – Review & Editing (supporting).

477 **Pietro Cozzini:** Software (equal), Writing – Review & Editing (supporting).

478 **Francesca Peccati:** Software (supporting), Writing – Review & Editing (supporting).

479 **Paolo P. Mazzeo:** Conceptualization (equal), Investigation (supporting), Data Curation (supporting), Validation (supporting),  
480 Writing – Review & Editing (supporting).

481 **Nicolò Riboni:** Investigation (supporting), Writing – Review & Editing (supporting).

482 **Maria Careri:** Supervision (equal), Funding Acquisition (equal), Resources (equal), Writing – Review & Editing (equal).

483 **Alessia Bacchi:** Conceptualization (lead), Funding Acquisition (equal), Resources (equal), Supervision (equal), Writing –  
484 Review & Editing (equal).

485

486 **Declaration of Competing Interest**

487 The authors declare that they have no known competing financial interests or personal relationships that could have appeared  
488 to influence the work reported in this paper.

489 **References**

- 490 [1] E.V.R. Campos, P.L.F. Proença, J.L. Oliveira, M. Bakshi, P.C. Abhilash, L.F. Fraceto, Use of botanical insecticides  
491 for sustainable agriculture: Future perspectives, *Ecol. Indic.* 105 (2019) 483–495.  
492 doi:10.1016/J.ECOLIND.2018.04.038.
- 493 [2] Y. Kourkoutas, M. Angane, S. Swift, K. Huang, C.A. Butts, S.Y. Quek, Essential Oils and Their Major  
494 Components: An Updated Review on Antimicrobial Activities, Mechanism of Action and Their Potential  
495 Application in the Food Industry, *Foods* 2022, Vol. 11, Page 464. 11 (2022) 464. doi:10.3390/FOODS11030464.
- 496 [3] N.S. Singh, R. Sharma, T. Parween, P.K. Patanjali, Pesticide Contamination and Human Health Risk Factor, *Mod.*  
497 *Age Environ. Probl. Their Remediat.* (2018) 49–68. doi:10.1007/978-3-319-64501-8\_3.
- 498 [4] *Climate Change, Intercropping, Pest Control and Beneficial Microorganisms*, Springer Netherlands, 2009.  
499 doi:10.1007/978-90-481-2716-0.
- 500 [5] F. Bakkali, S. Averbeck, D. Averbeck, M. Idaomar, Biological effects of essential oils - A review, *Food Chem.*  
501 *Toxicol.* 46 (2008) 446–475. doi:10.1016/j.fct.2007.09.106.
- 502 [6] M. Alonso-Gato, G. Astray, J.C. Mejuto, J. Simal-Gandara, Essential Oils as Antimicrobials in Crop Protection,  
503 *Antibiotics.* 10 (2021) 34. doi:10.3390/antibiotics10010034.
- 504 [7] F. Bianchi, F. Fornari, N. Riboni, C. Spadini, C.S. Cabassi, M. Iannarelli, C. Carraro, P.P. Mazzeo, A. Bacchi, S.  
505 Orlandini, S. Furlanetto, M. Careri, Development of novel cocrystal-based active food packaging by a Quality by  
506 Design approach, *Food Chem.* 347 (2021) 129051. doi:10.1016/j.foodchem.2021.129051.
- 507 [8] A.T.H. Mossa, Green Pesticides: Essential oils as biopesticides in insect-pest management, *J. Environ. Sci. Technol.*  
508 9 (2016) 354–378. doi:10.3923/jest.2016.354.378.
- 509 [9] S. Sharma, S. Barkauskaite, A.K. Jaiswal, S. Jaiswal, Essential oils as additives in active food packaging, *Food*  
510 *Chem.* (2020) 128403. doi:10.1016/j.foodchem.2020.128403.
- 511 [10] Food and Drug Administration, 54960 Federal Register / Vol . 81 , No . 159 / Wednesday , August 17 , 2016 / Rules  
512 and Regulations, 81 (2016) 54960–55055.
- 513 [11] J. Wiczyńska, I. Cavoski, Antimicrobial, antioxidant and sensory features of eugenol, carvacrol and trans-anethole  
514 in active packaging for organic ready-to-eat iceberg lettuce, *Food Chem.* 259 (2018) 251–260.  
515 doi:10.1016/j.foodchem.2018.03.137.
- 516 [12] L. Pavoni, D.R. Perinelli, G. Bonacucina, M. Cespi, G.F. Palmieri, An overview of micro- and nanoemulsions as

- 517 vehicles for essential oils: Formulation, preparation and stability, *Nanomaterials*. 10 (2020) 135.  
518 doi:10.3390/nano10010135.
- 519 [13] G.R. Desiraju, *Crystal Engineering: A Holistic View*, *Angew. Chemie Int. Ed.* 46 (2007) 8342–8356.  
520 doi:10.1002/anie.200700534.
- 521 [14] A. Bacchi, P.P. Mazzeo, *Cocrystallization as a tool to stabilize liquid active ingredients*, *Crystallogr. Rev.* (2021) 1–  
522 22. doi:10.1080/0889311X.2021.1978079.
- 523 [15] D. Balestri, P.P. Mazzeo, R. Perrone, F. Fornari, F. Bianchi, M. Careri, A. Bacchi, P. Pelagatti, *Deciphering the*  
524 *Supramolecular Organization of Multiple Guests Inside a Microporous MOF to Understand their Release Profile*,  
525 *Angew. Chemie Int. Ed.* 60 (2021) 10194–10202. doi:10.1002/ANIE.202017105.
- 526 [16] D. Balestri, P.P. Mazzeo, C. Carraro, N. Demitri, P. Pelagatti, A. Bacchi, *Stepwise Evolution of Molecular*  
527 *Nanoaggregates Inside the Pores of a Highly Flexible Metal–Organic Framework*, *Angew. Chemie*. 131 (2019)  
528 17503–17511. doi:10.1002/ange.201907621.
- 529 [17] P.P. Mazzeo, S. Canossa, C. Carraro, P. Pelagatti, A. Bacchi, *Systematic coformer contribution to cocrystal*  
530 *stabilization: energy and packing trends*, *CrystEngComm*. 22 (2020) 7341–7349. doi:10.1039/D0CE00291G.
- 531 [18] P.P. Mazzeo, C. Carraro, A. Arns, P. Pelagatti, A. Bacchi, *Diversity through Similarity: A World of Polymorphs,*  
532 *Solid Solutions, and Cocrystals in a Vial of 4,4'-Diazopyridine*, *Cryst. Growth Des.* 20 (2020) 636–644.  
533 doi:10.1021/acs.cgd.9b01052.
- 534 [19] P.P. Mazzeo, C. Carraro, A. Monica, D. Capucci, P. Pelagatti, F. Bianchi, S. Agazzi, M. Careri, A. Raio, M. Carta,  
535 F. Menicucci, M. Belli, M. Michelozzi, A. Bacchi, *Designing a Palette of Cocrystals Based on Essential Oil*  
536 *Constituents for Agricultural Applications*, *ACS Sustain. Chem. Eng.* 7 (2019) 17929–17940.  
537 doi:10.1021/acssuschemeng.9b04576.
- 538 [20] D. Capucci, D. Balestri, P.P. Mazzeo, P. Pelagatti, K. Rubini, A. Bacchi, *Liquid Nicotine Tamed in Solid Forms by*  
539 *Cocrystallization*, *Cryst. Growth Des.* 17 (2017) 4958–4964. doi:10.1021/acs.cgd.7b00887.
- 540 [21] P.P. Mazzeo, M. Pioli, F. Montisci, A. Bacchi, P. Pelagatti, *Mechanochemical Preparation of Dipyridyl-*  
541 *Naphthalenediimide Cocrystals: Relative Role of Halogen-Bond and  $\pi$ - $\pi$  Interactions*, *Cryst. Growth Des.* 21  
542 (2021) 5687–5696. doi:10.1021/acs.cgd.1c00531.
- 543 [22] N.K. Duggirala, M.L. Perry, Ö. Almarsson, M.J. Zaworotko, *Pharmaceutical cocrystals: Along the path to improved*  
544 *medicines*, *Chem. Commun.* 52 (2016) 640–655. doi:10.1039/c5cc08216a.

- 545 [23] J.W. Steed, The role of co-crystals in pharmaceutical design., *Trends Pharmacol. Sci.* 34 (2013) 185–93.  
546 doi:10.1016/j.tips.2012.12.003.
- 547 [24] J.G.P. Wicker, L.M. Crowley, O. Robshaw, E.J. Little, S.P. Stokes, R.I. Cooper, S.E. Lawrence, Will they co-  
548 crystallize?, *CrystEngComm*. 19 (2017) 5336–5340. doi:10.1039/C7CE00587C.
- 549 [25] A. Bacchi, D. Capucci, M. Giannetto, M. Mattarozzi, P. Pelagatti, N. Rodriguez-Hornedo, K. Rubini, A. Sala,  
550 Turning Liquid Propofol into Solid (without Freezing It): Thermodynamic Characterization of Pharmaceutical  
551 Cocrystals Built with a Liquid Drug, *Cryst. Growth Des.* 16 (2016) 6547–6555. doi:10.1021/acs.cgd.6b01241.
- 552 [26] Y. Xiao, L. Zhou, H. Hao, Y. Bao, Q. Yin, C. Xie, Cocrystals of propylthiouracil and nutraceuticals toward  
553 sustained-release: Design, structure analysis, and solid-state characterization, *Cryst. Growth Des.* 21 (2021) 1202–  
554 1217. doi:10.1021/acs.cgd.0c01519.
- 555 [27] O. Shemchuk, S. D’Agostino, C. Fiore, V. Sambri, S. Zannoli, F. Grepioni, D. Braga, Natural Antimicrobials Meet  
556 a Synthetic Antibiotic: Carvacrol/Thymol and Ciprofloxacin Cocrystals as a Promising Solid-State Route to  
557 Activity Enhancement, *Cryst. Growth Des.* 20 (2020) 6796–6803. doi:10.1021/acs.cgd.0c00900.
- 558 [28] M.D. Perera, J. Desper, A.S. Sinha, C.B. Aakeröy, Impact and importance of electrostatic potential calculations for  
559 predicting structural patterns of hydrogen and halogen bonding, *CrystEngComm*. 18 (2016) 8631–8636.  
560 doi:10.1039/c6ce02089e.
- 561 [29] M.C. Etter, Encoding and Decoding Hydrogen-Bond Patterns of Organic Compounds, *Acc. Chem. Res.* 23 (1990)  
562 120–126. doi:10.1021/ar00172a005.
- 563 [30] C.A. Hunter, Quantifying intermolecular interactions: Guidelines for the molecular recognition toolbox, *Angew.*  
564 *Chemie - Int. Ed.* 43 (2004) 5310–5324. doi:10.1002/anie.200301739.
- 565 [31] M. Karimi-Jafari, L. Padrela, G.M. Walker, D.M. Croker, Creating cocrystals: A review of pharmaceutical cocrystal  
566 preparation routes and applications, *Cryst. Growth Des.* 18 (2018) 6370–6387. doi:10.1021/acs.cgd.8b00933.
- 567 [32] N. Issa, P.G. Karamertzanis, G.W.A. Welch, S.L. Price, Can the formation of pharmaceutical cocrystals be  
568 computationally predicted? I. Comparison of lattice energies, *Cryst. Growth Des.* 9 (2009) 442–453.  
569 doi:10.1021/cg800685z.
- 570 [33] M.A. Mohammad, A. Alhalaweh, S.P. Velaga, Hansen solubility parameter as a tool to predict cocrystal formation,  
571 *Int. J. Pharm.* 407 (2011) 63–71. doi:10.1016/j.ijpharm.2011.01.030.
- 572 [34] L. Fábíán, Cambridge Structural Database Analysis of Molecular Complementarity in Cocrystals, *Cryst. Growth*

573 Des. 9 (2009) 1436–1443. doi:10.1021/cg800861m.

574 [35] C.R. Groom, I.J. Bruno, M.P. Lightfoot, S.C. Ward, The Cambridge structural database, *Acta Crystallogr. Sect. B*  
575 *Struct. Sci. Cryst. Eng. Mater.* 72 (2016) 171–179. doi:10.1107/S2052520616003954.

576 [36] J. Devogelaer, H. Meekes, P. Tinnemans, E. Vlieg, R. Gelder, Co- crystal Prediction by Artificial Neural  
577 Networks\*\*, *Angew. Chemie Int. Ed.* 59 (2020) 21711–21718. doi:10.1002/anie.202009467.

578 [37] D. Wang, Z. Yang, B. Zhu, X. Mei, X. Luo, Machine-Learning-Guided Cocrystal Prediction Based on Large Data  
579 Base, *Cryst. Growth Des.* 20 (2020) 6610–6621. doi:10.1021/acs.cgd.0c00767.

580 [38] M. Przybyłek, T. Jeliński, J. Słabuszewska, D. Ziółkowska, K. Mroczyńska, P. Cysewski, Application of  
581 Multivariate Adaptive Regression Splines (MARSplines) Methodology for Screening of Dicarboxylic Acid  
582 Cocrystal Using 1D and 2D Molecular Descriptors, *Cryst. Growth Des.* 19 (2019) 3876–3887.  
583 doi:10.1021/acs.cgd.9b00318.

584 [39] M. Przybyłek, P. Cysewski, Distinguishing Cocrystals from Simple Eutectic Mixtures: Phenolic Acids as Potential  
585 Pharmaceutical Coformers, *Cryst. Growth Des.* 18 (2018) 3524–3534. doi:10.1021/acs.cgd.8b00335.

586 [40] A. Vriza, A.B. Canaj, R. Vismara, L.J. Kershaw Cook, T.D. Manning, M.W. Gaultois, P.A. Wood, V. Kurlin, N.  
587 Berry, M.S. Dyer, M.J. Rosseinsky, One class classification as a practical approach for accelerating  $\pi$ - $\pi$  co-crystal  
588 discovery, *Chem. Sci.* 12 (2021) 1702–1719. doi:10.1039/d0sc04263c.

589 [41] M.E. Mswahili, M.J. Lee, G.L. Martin, J. Kim, P. Kim, G.J. Choi, Y.S. Jeong, Cocrystal prediction using machine  
590 learning models and descriptors, *Appl. Sci.* 11 (2021) 1–12. doi:10.3390/app11031323.

591 [42] H. Moriwaki, Y.S. Tian, N. Kawashita, T. Takagi, Mordred: A molecular descriptor calculator, *J. Cheminform.* 10  
592 (2018) 4. doi:10.1186/s13321-018-0258-y.

593 [43] J.G.P. Wicker, R.I. Cooper, Will it crystallise? Predicting crystallinity of molecular materials, *CrystEngComm.* 17  
594 (2015) 1927–1934. doi:10.1039/c4ce01912a.

595 [44] I.J. Bruno, J.C. Cole, P.R. Edgington, M. Kessler, C.F. Macrae, P. McCabe, J. Pearson, R. Taylor, New software for  
596 searching the Cambridge Structural Database and visualizing crystal structures, *Acta Crystallogr. Sect. B Struct.*  
597 *Sci.* 58 (2002) 389–397. doi:10.1107/S0108768102003324.

598 [45] C.F. Macrae, I.J. Bruno, J.A. Chisholm, P.R. Edgington, P. McCabe, E. Pidcock, L. Rodriguez-Monge, R. Taylor, J.  
599 van de Streek, P.A. Wood, Mercury CSD 2.0 – new features for the visualization and investigation of crystal  
600 structures, *J. Appl. Crystallogr.* 41 (2008) 466–470. doi:10.1107/S0021889807067908.



- 601 [46] M. Baroni, G. Cruciani, S. Sciabola, F. Perruccio, J.S. Mason, A Common Reference Framework for  
602 Analyzing/Comparing Proteins and Ligands. Fingerprints for Ligands And Proteins (FLAP): Theory and  
603 Application, *J. Chem. Inf. Model.* 47 (2007) 279–294. doi:10.1021/ci600253e.
- 604 [47] Sybyl 8.1, Tripos International, St. Louis, USA, 2009.
- 605 [48] M. Clark, R.D. Cramer, N. Van Opdenbosch, Validation of the general purpose tripos 5.2 force field, *J. Comput.*  
606 *Chem.* 10 (1989) 982–1012. doi:10.1002/jcc.540100804.
- 607 [49] Schrödinger LLC, The PyMOL Molecular Graphics System, PyMOL Mol. Graph. Syst. Version 2.0. (2010).
- 608 [50] G. Landrum, RDKit: Open-Source Cheminformatics Software, [Http://Www.Rdkit.Org/](http://www.rdkit.org/). (2021).
- 609 [51] G.A. Frisch, M. J., Trucks, G. W., Schlegel, H. B., Scuseria, G. E., Robb, M. A., Cheeseman, J. R., Scalmani, G.,  
610 Barone, V., Petersson, Gaussian 16, Rev. A.03, Gaussian, Inc., Wallingford, CT. (2016).
- 611 [52] R.D. Snee, Validation of Regression Models: Methods and Examples, *Technometrics.* 19 (1977) 415–428.  
612 doi:10.1080/00401706.1977.10489581.
- 613 [53] M. Li Vigni, C. Durante, M. Cocchi, Exploratory Data Analysis, in: *Data Handl. Sci. Technol.*, Elsevier Ltd, 2013:  
614 pp. 55–126. doi:10.1016/B978-0-444-59528-7.00003-X.
- 615 [54] S. Wold, K. Esbensen, P. Geladi, Principal component analysis, *Chemom. Intell. Lab. Syst.* 2 (1987) 37–52.  
616 doi:10.1016/0169-7439(87)80084-9.
- 617 [55] R. Bro, A.K. Smilde, Principal component analysis, *Anal. Methods.* 6 (2014) 2812–2831. doi:10.1039/c3ay41907j.
- 618 [56] M. Cocchi, A. Biancolillo, F. Marini, Chemometric Methods for Classification and Feature Selection, in: *Compr.*  
619 *Anal. Chem.*, Elsevier B.V., 2018: pp. 265–299. doi:10.1016/bs.coac.2018.08.006.
- 620 [57] U.G. Indahl, H. Martens, T. Næs, From dummy regression to prior probabilities in PLS-DA, *J. Chemom.* 21 (2007)  
621 529–536. doi:10.1002/cem.1061.
- 622 [58] S. Wold, M. Sjöström, L. Eriksson, PLS-regression: A basic tool of chemometrics, in: *Chemom. Intell. Lab. Syst.*,  
623 Elsevier, 2001: pp. 109–130. doi:10.1016/S0169-7439(01)00155-1.
- 624 [59] D. Braga, E. Dichiarante, F. Grepioni, G.I. Lampronti, L. Maini, P.P. Mazzeo, S. D’Agostino, Mechanical  
625 Preparation of Crystalline Materials. An Oxymoron?, in: *Supramol. Chem.*, John Wiley & Sons, Ltd, Chichester,  
626 UK, 2012. doi:10.1002/9780470661345.smc115.
- 627 [60] J. Fürtkranz, P.K. Chan, S. Craw, C. Sammut, W. Uther, A. Ratnaparkhi, X. Jin, J. Han, Y. Yang, K. Morik, M.  
628 Dorigo, M. Birattari, T. Stütze, P. Brazdil, R. Vilalta, C. Giraud-Carrier, C. Soares, J. Rissanen, R.A. Baxter, I.

629 Bruha, R.A. Baxter, G.I. Webb, L. Torgo, A. Banerjee, H. Shan, S. Ray, P. Tadepalli, Y. Shoham, R. Powers, Y.  
630 Shoham, R. Powers, G.I. Webb, S. Ray, S. Scott, H. Blockeel, L. De Raedt, Manhattan Distance, *Encycl. Mach.*  
631 *Learn.* (2011) 639–639. doi:10.1007/978-0-387-30164-8\_506.

632 [61] L.E. Juarez-Orozco, O. Martinez-Manzanera, S. V. Nesterov, S. Kajander, J. Knuuti, The machine learning horizon  
633 in cardiac hybrid imaging, *Eur. J. Hybrid Imaging.* 2 (2018) 1–15. doi:10.1186/S41824-018-0033-3/FIGURES/5.

634 [62] A. Singh, N. Thakur, A. Sharma, A review of supervised machine learning algorithms, in: *3rd Int. Conf. Comput.*  
635 *Sustain. Glob. Dev.*, Institute of Electrical and Electronics Engineers, New Delhi, India, 2016: pp. 1310–1315.  
636 <https://ieeexplore.ieee.org/abstract/document/7724478>.

637 [63] N.M. Faber, R. Bro, Standard error of prediction for multiway PLS: 1. Background and a simulation study,  
638 *Chemom. Intell. Lab. Syst.* 61 (2002) 133–149. doi:10.1016/S0169-7439(01)00204-0.

639