

Supplementary Material:

Analogues of the Frog-skin Antimicrobial Peptide Temporin 1Tb Exhibit a Wider Spectrum of Activity and a Stronger Antibiofilm Potential as Compared to the Parental Peptide

Lucia Grassi, Giuseppantonio Maisetta, Giuseppe Maccari, Semih Esin, Giovanna Batoni *

***Correspondence:**

Giovanna Batoni
giovanna.batoni@med.unipi.it

1 SUPPLEMENTARY DATA

1.1 Dataset preparation

A dataset representing peptides with cytotoxic activity was designed with the aim to train and validate a statistical model able to discern between ‘toxic’ and ‘non-toxic’ peptides, giving a confidence score. A set of sequences ranging from 9 to 35 amino acids length was collected from different bioactive peptide databases, as previously described (Gupta et al, 2013). After removal of peptides with non-standard residues, 1709 peptides were left. The negative dataset was populated with non-secretory sequences randomly extracted from UniProt database, without the ‘antimicrobial’ and ‘cytotoxic’ annotation and with a length ranging from 9 to 35 amino acids, for a total count of 2010 negative sequences. A homology cut-off was imposed to exclude similar peptides in order to avoid redundant data that could influence the prediction performance. Peptides showing a sequence identity equal or greater than 70% to any other in the dataset were identified and removed by the CD-HIT (Cluster Database at High Identity with Tolerance) program (Li and Godzik, 2006).

1.2 Data encoding

In order to build a statistical model, able to discern between toxic and non-toxic peptides, each sequence in the dataset was encoded into computer-intelligible variables representing peptides physicochemical peculiarities. Peptide charge at different pH conditions, isoelectric point and molecular weight, together with the z-scale moment, were used to describe global features of the peptide sequences. Z-scale descriptors (Hellberg et al, 1987) are highly condensed variables, originally derived from a principal component analysis (PCA) of several experimental and theoretical physicochemical properties for the 20 naturally occurring amino acids (AAs). These descriptors were successively expanded to include artificial AAs for a total of 87 AAs (Sandberg et al, 1998). In detail, this latter version corresponds to the first five principal components explaining the variance in the set: z_1 , z_2 , and z_3 represent the AA hydrophobicity, steric properties, and polarity, respectively, while z_4 and z_5 describe the electronic effects of the residues. The z-scale moment (μZ_i), an extension of Eisenberg’s hydrophobic moment equation (Eisenberg et al, 1982), represents z-scales distribution along peptide sequences.

$$\mu Z_i = \sqrt{\left(\sum_{k=1}^L Z_i^k \sin(\delta k) \right)^2 + \left(\sum_{k=1}^L Z_i^k \cos(\delta k) \right)^2}$$

Equation 1. Z-scale moment

In Equation 1, δ is the angular frequency of the AA residues forming the structure (100° for alpha helix); k is the number of the particular residue examined, L is the length of the sequence and Z_i^k is the z_i -scale value of the k^{th} AA. In particular, μZ_1 represents a measure of the hydrophobicity distribution along peptide sequence. Topological descriptors represent the interaction of different residues along the amino acidic sequence and are used to keep into account peptide's secondary structure. QSAR descriptors were encoded into auto- and cross covariance (ACC) values. Classical ACC transformation was introduced by Wold et al. (Wold et al, 1993) and results in two kinds of variables: auto covariance (AC) of the same descriptor and cross covariance (CC) between two different descriptors. Briefly, for a given protein sequence, ACC variables describe the average interactions between residues distributed a certain *lag* apart throughout the whole sequence. In this work, the *Minimum and Maximum of auto- and cross-covariances* (mMACC) algorithm is used (Maccari et al, 2013), weak and strong correlations are kept into account (Equation 2).

$$AC_{\min d} = MIN[Z_i^k * Z_i^{k+d}] \quad AC_{\max d} = MAX[Z_i^k * Z_i^{k+d}] \quad (k = 1, 2, 3..L-d)$$

$$CC_{\min d} = MIN[Z_i^k * Z_i^{k+d}] \quad CC_{\max d} = MAX[Z_i^k * Z_i^{k+d}] \quad (k = 1, 2, 3..L-d)$$

Equation 2. Minimum and Maximum of auto and cross-covariance equations

Both in the global and topological descriptors, Z-scale values were mean-centered and scaled prior to their use, as described by the following equation:

$$Z_i = \frac{z_i - \frac{1}{N} \sum_{k=1}^N z_i^k}{\sqrt{\frac{1}{N} \sum_{j=1}^N \left[z_i^j - \frac{1}{N} \sum_{k=1}^N z_i^k \right]^2}}$$

Equation 3. Z-scale descriptor normalization

Where Z_i is the i^{th} descriptor of z-scales variables, z_i is the original z-scale value and N is the number of AAs in the z-scales descriptors table.

1.3 Feature selection and model generation

In this study, the Random Forest algorithm (RF), implemented in the software suite WEKA (Witten et al, 2011), was adopted as prediction engine. Model performance was measured with a 10-fold cross-validation analysis, where each dataset was divided into 10 parts - 9 parts for model learning (training) and the remaining part for validation (testing). As a performance measure, the Matthews correlation coefficient (MCC) was used, as defined below.

$$\text{Sensitivity} = \frac{TP}{TP + FN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$MCC = \frac{(TP * TN) - (FN * FP)}{\sqrt{(TP + FN) * (TN + FP) * (TP + FP) * (TN + FN)}}$$

Equation 4. Performance evaluation equations

Where TP , TN , FP and FN are the number of true positive, true negative, false positive and false negative, respectively, resulting from the model. MCC is an important index used to evaluate the performance of the predictor when the dataset is not balanced (Baldi et al, 2000). In order to obtain a non-redundant set of descriptors, the Maximum Relevance, Minimum Redundancy (mRMR) method (Peng et al, 2005) was employed to sort features in descending order of importance. Incremental Feature Selection (IFS) (Huang et al, 2010) was applied to the sorted descriptors list by consecutively incrementing by 5 the number of descriptors. Each descriptor set thus obtained was evaluated by tenfold cross-validation and the IFS curve was plotted to unveil the relation between the performance of the model and the feature subset. The optimal feature subset is defined as that showing the highest MCC value (**Figure S1**); the selected model was used for peptides classification. A description of the applied descriptors is available in **Table S1**, while the hierarchical list of the final descriptors is shown in **Table S2**.

1.4 Sequence similarity

For TB peptide optimization, a supplemental objective representing sequence similarity was added. Sequence similarity is defined by the Smith-Waterman score between the respective peptide sequences (Smith at al, 1981). However, since the Smith-Waterman score is dependent on input sequences length, the final score was normalized between 0 and 1 by dividing by the maximum score of the two self-alignments, as shown in Equation 5 (Zang et al, 2012).

$$NS_{A,B} = \frac{S_{A,B}}{\max(S_{A,A}, S_{B,B})}$$

Equation 5. Smith-Waterman normalized score

Here, $S_{A,B}$ is the similarity score between sequence A and B, $S_{A,A}$ and $S_{B,B}$ are the self-alignment score of sequence A and sequence B, respectively. In order to consider not only the identity between two amino acidic positions, a score matrix was defined by calculating the Euclidean distance between the five auto-scaled z-scale values of each AA pairs.

References

- Baldi P., Brunak S., Chauvin Y., Andersen C.A.F., Nielsen H. (2000). Assessing the accuracy of prediction algorithms for classification: an overview. *Bioinformatics*. 16: 412-424.
- Eisenberg D., Weiss R.M., Terwilliger T.C. (1982). The helical hydrophobic moment: a measure of the amphiphilicity of a helix. *Nature* 299: 371-374.
- Gupta, S., Kapoor, P., Chaudhary, K., Gautam, A., Kumar, R., Raghava, G.P.S. (2013). In silico approach for predicting toxicity of peptides and proteins. *PLoS One*. 8:e73957. doi: 10.1371/journal.pone.0073957.
- Hellberg S., Sjöström M., Skagerberg B., Wold S. (1987). Peptide quantitative structure activity relationship, a multivariate approach. *J. Med. Chem.* 30: 1126-1135.
- Huang T., Shi X.H., Wang P., He Z., Feng K.Y., Hu L., Kong X., Li Y.X., Cai Y.D., Chou K.C. (2010). Analysis and prediction of the metabolic stability of proteins based on their sequential features, subcellular locations and interaction networks. *PLoS One*. 5: e10972.
- Li, W., Godzik, A. (2006). Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide sequences. *Bioinformatics*. 22: 1658-1659. doi: 10.1093/bioinformatics/btl158.
- Maccari G., Di Luca M., Nifosí R., Cardarelli F., Signore G., Boccardi C., Bifone A. (2013). Antimicrobial peptides design by evolutionary multiobjective optimization. *PLoS Comput Biol*. 9: e1003212.
- Peng H., Long F., Ding C. (2005). Feature selection based on mutual information: criteria of max-dependency, max-relevance, and min-redundancy. *IEEE Trans. Pattern Anal. Mach. Intell.* 27: 1226-1238.
- Sandberg M., Eriksson L., Jonsson J., Sjöström M., Wold S. (1998). New chemical descriptors relevant for the design of biologically active peptides. A multivariate characterization of 87 amino acids. *J. Med. Chem.* 41: 2481-2491.
- Smith T.F., Waterman M.S. (1981). Identification of common molecular subsequences. *J. Mol. Biol.* 147: 195–197.
- Witten I.H., Frank E., Hall M.A. (2011). Data mining: practical machine learning tools and techniques. Morgan Kaufmann, Burlington, MA
- Wold S., Jonsson J., Sjöström M., Sandberg M., Rännar S. (1993). DNA and peptide sequences and chemical processes multivariately modelled by principal component analysis and partial least-squares projections to latent structures. *Ann. Chim. Acta*. 277: 239-253.
- Zang M., Leong H. (2012). BBH-LS: an algorithm for computing positional homologs using sequence and gene context similarity. *BMC Syst. Biol.* 6: S22.

2 SUPPLEMENTARY TABLES AND FIGURES

2.1 Tables

Table S1. Applied descriptors in the model building. A list of the applied descriptors with abbreviation and description is provided.

Type	Abbreviation	Description
Global	NetCharge@5	Net charge at pH = 5.
	NetCharge@7	Net charge at pH = 7.
	NetCharge@9	Net charge at pH = 9.
	Wimley White (pH n)	Wimley White partitioning at pH <i>n</i>
	Isoelectric point	Peptide's isoelectric point
	Size	Total amino acid count.
	Property_Zn	Z-scale average sum of property <i>n</i> along peptide sequence.
	Variable Moment Zn (100 D)	Z-scale moment distribution of property <i>n</i> along peptide sequence at 100 degrees (the angle between two residues in alpha helix conformation)
	D_X_AC_LAG_N_[MIN,MAX]; D_X:Y_CC_LAG_N_[MIN,MAX]	D_X_AC_LAG_N_[MIN,MAX]: Topological descriptor of the auto covariance of descriptor X with a lag of N. D_X:Y_CC_LAG_N_[MIN,MAX]: Topological descriptor of the cross covariance between descriptor X and Y, with a lag of N. With X and Y being a value between 0 and 4: 0) Z-scale Descriptor 1 1) Z-scale Descriptor 2 2) Z-scale Descriptor 3 3) Z-scale Descriptor 4 4) Z-scale Descriptor 5

Table S2. Hierarchical list of descriptors. List of descriptors sorted by the mRMR method.

Name	#	Name	#	Name	#	Name	#	Name	#	Name
2 Property_z5_10	11 D_4:0_CC_LAG_3_	23 7 MIN	23 11 D_0:1_CC_LAG_4_	23 8 MIN	23 3 D_0:4_CC_LAG_5_MIN	34 7 D_1:3_CC_LAG_5_MIN	46 2 D_1:0_CC_LAG_7_MAX	57 57 D_4:3_CC_LAG_1_-	7 MIN	
3 Property_z2_6	11 8 MIN	23 3 D_0:4_CC_LAG_5_MIN	34 8 D_1_AC_LAG_4_MIN	46 3 D_2:1_CC_LAG_3_MAX	46 46 D_2:1_CC_LAG_3_MAX	57 8 MAX	57 D_3:1_CC_LAG_9_-	8 MAX	57	
4 Property_z3_14	9 Property_z2_12	12 4 D_0:4_CC_LAG_8_MAX	23 4 D_0:4_CC_LAG_8_MAX	35 9 D_2:0_CC_LAG_7_MIN	46 4 D_0_AC_LAG_8_MIN	58 9 Property_z3_20	58 58 D_2:3_CC_LAG_0_-	9 Property_z3_20	7 MIN	
5 Property_z2_2	0 Property_z3_16	12 5 D_1:0_CC_LAG_4_MIN	23 0 D_0_AC_LAG_6_MIN	35 0 D_0_AC_LAG_6_MIN	46 5 D_3:1_CC_LAG_1_MAX	58 0 MAX	58 58 D_2:3_CC_LAG_5_-	0 MAX	58	
6 Property_z5_9	1 N	12 6 D_2:1_CC_LAG_3_MIN	23 1 D_2:1_CC_LAG_6_MAX	35 1 D_2:1_CC_LAG_6_MAX	46 6 D_4:3_CC_LAG_8_MIN	58 1 MIN	58 58 D_4:3_CC_LAG_4_-	1 MIN	58	
7 Property_z2_3	12 D_0:2_CC_LAG_0_-	23 23 MIN	7 7 D_4_AC_LAG_6_MAX	35 2 D_0:2_CC_LAG_8_MIN	46 7 Property_z1_31	58 2 Property_z1_29	58 58 D_1:2_CC_LAG_2_-	2 Property_z1_29	7 MIN	
8 Property_z2_1	12 12 D_1:0_CC_LAG_1_-	23 3 MIN	8 8 D_3:0_CC_LAG_3_MAX	35 3 D_4:0_CC_LAG_9_MIN	46 8 D_3_AC_LAG_4_MAX	58 3 MAX	58 58 D_2:1_CC_LAG_9_-	3 MAX	58	
9 Property_z5_6	12 4 MAX	23 9 D_0:3_CC_LAG_4_MIN	23 4 D_1:0_CC_LAG_2_MAX	35 9 D_4:2_CC_LAG_9_MAX	46 9 D_4:2_CC_LAG_9_MAX	58 4 MAX	58 58 D_4:3_CC_LAG_4_-	4 MAX	58	
10 Property_z4_5	5 Property_z4_34	12 12 D_0:4_CC_LAG_1_-	24 0 D_1:4_CC_LAG_3_MIN	35 5 D_2:0_CC_LAG_3_MAX	47 0 D_1_AC_LAG_5_MAX	58 5 MIN	58 58 D_2:4_CC_LAG_8_-	5 MIN	58	
11 Property_z2_34	6 MAX	12 1 D_2:1_CC_LAG_2_-	24 1 D_1:3_CC_LAG_3_MAX	35 6 D_3:1_CC_LAG_9_MIN	47 1 D_3:2_CC_LAG_5_MIN	58 6 MIN	58 58 D_4:2_CC_LAG_8_-	6 MIN	58	
12 Property_z2_9	7 MIN	12 2 D_0:4_CC_LAG_1_-	24 2 D_4_AC_LAG_4_MAX	35 7 D_3_AC_LAG_0_MIN	47 2 Property_z2_16	58 7 Property_z2_18	58 58 D_4:2_CC_LAG_0_-	7 Property_z2_18	7 MIN	
13 Property_z2_7	8 MIN	12 3 D_3:0_CC_LAG_1_-	24 3 Property_z5_14	35 8 D_0_AC_LAG_7_MAX	47 3 D_2:1_CC_LAG_9_MIN	58 8 MIN	58 58 D_3:2_CC_LAG_3_-	8 MIN	58	
14 Property_z1_4	9 MAX	12 4 D_0:2_CC_LAG_8_MAX	24 4 D_0:2_CC_LAG_8_MAX	36 9 D_1:4_CC_LAG_5_MIN	47 4 D_1_AC_LAG_6_MAX	59 9 MAX	59 59 D_4:2_CC_LAG_5_-	9 MAX	59	
15 Property_z2_11	0 Property_z5_12	13 5 D_4:0_CC_LAG_7_-	24 5 Property_z5_33	36 0 Property_z5_32	47 5 D_3:4_CC_LAG_2_MIN	0 Property_z2_29	59 59 D_4:2_CC_LAG_5_-	0 Property_z2_29	59	
16 Property_z3_1	1 MAX	13 6 D_0:1_CC_LAG_3_MAX	24 6 D_0:1_CC_LAG_3_MAX	36 1 D_2:0_CC_LAG_6_MIN	47 6 D_3_AC_LAG_8_MAX	1 Property_z3_22	59 59 D_0:2_CC_LAG_7_-	1 Property_z3_22	59	
17 N	2 Property_z4_0	13 7 D_0:1_CC_LAG_0_-	24 7 D_0:1_CC_LAG_7_MIN	36 2 D_1:2_CC_LAG_0_MAX	47 7 D_1:2_CC_LAG_8_MAX	2 MAX	59 59 D_4:2_CC_LAG_5_-	2 MAX	59	
18 Property_z1_9	3 MIN	13 8 D_4:1_CC_LAG_1_-	24 8 D_0_AC_LAG_1_MIN	36 3 D_3:2_CC_LAG_0_MIN	47 8 Property_z2_31	3 MIN	59 59 D_2:4_CC_LAG_5_-	3 MIN	59	
19 Property_z3_3	4 MAX	13 9 D_4:3_CC_LAG_7_MAX	25 9 D_4:3_CC_LAG_7_MAX	36 4 D_1:4_CC_LAG_9_MAX	48 9 D_1:3_CC_LAG_8_MAX	4 MIN	59 59 D_2:4_CC_LAG_5_-	4 MIN	59	
20 Property_z2_10	5 MIN	13 10 D_4:1_CC_LAG_1_-	25 10 D_3:1_CC_LAG_2_MIN	36 5 D_1:3_CC_LAG_3_MIN	48 0 D_2:1_CC_LAG_8_MAX	5 Property_z3_21	59 59 D_3:4_CC_LAG_5_-	5 Property_z3_21	59	
21 Property_z1_2	6 Property_z2_13	13 11 D_3:0_CC_LAG_4_-	25 11 D_0:2_CC_LAG_4_MIN	36 6 D_0:3_CC_LAG_8_MIN	48 1 D_2:3_CC_LAG_9_MIN	6 MIN	59 59 D_2:0_CC_LAG_0_-	6 MIN	59	
22 Property_z2_8	7 MIN	13 12 D_3_AC_LAG_7_MIN	25 12 D_3_AC_LAG_7_MIN	36 7 D_0:1_CC_LAG_8_MAX	48 2 D_4:2_CC_LAG_5_MAX	7 MAX	59 59 D_2:0_CC_LAG_0_-	7 MAX	59	

47	Property_z1_1	16 D_2:0_CC_LAG_1_ 2 MIN 16 D_1:0_CC_LAG_8_	27 7 D_1:3_CC_LAG_7_MIN 27 8 D_4_AC_LAG_8_MIN 27 9 D_1:0_CC_LAG_0_MAX	39 2 D_3:1_CC_LAG_3_MIN 39 3 D_1:3_CC_LAG_5_MAX 39 4 D_1:2_CC_LAG_9_MIN 39 5 D_0:3_CC_LAG_8_MAX 39 6 D_4:2_CC_LAG_6_MAX 39 7 D_3:1_CC_LAG_5_MAX 39 8 D_0:2_CC_LAG_9_MIN 39 9 D_3:0_CC_LAG_7_MAX 40 0 D_4:1_CC_LAG_8_MIN 40 1 D_3:1_CC_LAG_2_MAX 40 2 D_2:4_CC_LAG_9_MAX 40 3 D_0_AC_LAG_8_MAX 40 4 D_2:4_CC_LAG_9_MIN 40 5 D_1:0_CC_LAG_4_MAX 40 6 D_3:4_CC_LAG_9_MIN 40 7 D_1:3_CC_LAG_0_MAX 40 8 D_4:2_CC_LAG_9_MIN 40 9 D_2:3_CC_LAG_1_MIN 41 0 D_0:3_CC_LAG_9_MIN 41 1 D_1_AC_LAG_7_MIN 41 2 D_1:2_CC_LAG_8_MIN 41 3 D_3:1_CC_LAG_0_MAX 41 4 D_2:3_CC_LAG_8_MIN 41 5 D_0:3_CC_LAG_9_MAX 41 6 D_3_AC_LAG_2_MAX	50 7 Property_z5_30 50 8 Property_z1_20 50 9 D_3:2_CC_LAG_8_MIN 51 0 D_2:3_CC_LAG_8_MAX 51 1 D_1_AC_LAG_6_MIN 51 2 D_2:4_CC_LAG_6_MAX 51 3 D_3_AC_LAG_3_MIN 51 4 Property_z1_19 51 5 D_0:2_CC_LAG_0_MAX 51 6 D_2:3_CC_LAG_0_MIN 51 7 D_1:3_CC_LAG_9_MAX 51 8 Property_z2_30 51 9 D_3_AC_LAG_8_MIN 52 0 D_3:4_CC_LAG_9_MIN 52 1 D_4:2_CC_LAG_9_MIN 52 2 Property_z3_18 52 3 D_2:4_CC_LAG_2_MAX 52 4 D_4:2_CC_LAG_1_MIN 52 5 D_2:0_CC_LAG_9_MAX 52 6 D_4:2_CC_LAG_5_	62 2 Property_z4_17 62 3 Property_z2_24 62 4 D_3:4_CC_LAG_4_ MI 62 5 D_0:2_CC_LAG_1_ MAX 62 6 Property_z2_28 62 7 D_4:2_CC_LAG_8_ MAX 62 8 D_2:4_CC_LAG_3_ MIN 62 9 D_2:0_CC_LAG_1_ MAX 63 0 Property_z2_21 63 1 D_3:4_CC_LAG_0_ MIN 63 2 Property_z4_28 63 3 D_2:3_CC_LAG_2_ MAX 63 4 Property_z2_23 63 5 D_3:2_CC_LAG_9_ MAX 63 6 D_0:2_CC_LAG_5_ MAX 63 7 D_2:3_CC_LAG_1_ MAX 63 8 Property_z2_22 63 9 D_3:4_CC_LAG_1_ MIN 64 0 Property_z3_24 64 1 D_2:3_CC_LAG_3_ MAX 64 2 Property_z4_27 64 3 D_2_AC_LAG_9_ MI 64 4 N 64 5 D_2:3_CC_LAG_2_ MAX 64 6 D_1_AC_LAG_9_MAX 64 7 D_3_AC_LAG_2_MAX
48	Property_z3_8	3 MIN	8 D_4_AC_LAG_8_MIN	39 3 D_1:3_CC_LAG_5_MAX	50 8 Property_z1_20	
49	D_3:1_CC_LAG_4_MI	16 D_4_AC_LAG_8_M	9 D_1:0_CC_LAG_0_MAX	50 9 D_3:2_CC_LAG_8_MIN	62 3 Property_z2_24	
50	N	4 AX	9 D_1:0_CC_LAG_0_MAX	51 0 D_2:3_CC_LAG_8_MAX	62 4 MIN	
51	D_4:3_CC_LAG_4_MA	16 D_3:1_CC_LAG_1_	0 D_0:4_CC_LAG_6_MIN	51 1 D_2:3_CC_LAG_8_MAX	62 5 MAX	
52	X	5 MIN	1 D_4_AC_LAG_2_MIN	51 2 D_1_AC_LAG_6_MIN	62 6 Property_z2_28	
53	D_3:0_CC_LAG_0_MI	16 D_1_AC_LAG_2_MI	1 D_4_AC_LAG_2_MIN	51 3 D_2:4_CC_LAG_6_MAX	62 7 MAX	
54	N	6 N	2 D_2:1_CC_LAG_6_MIN	51 4 D_3_AC_LAG_3_MIN	62 8 D_2:4_CC_LAG_3_ MIN	
55	Variable moment	16 D_4:0_CC_LAG_6_	2 D_2:1_CC_LAG_6_MIN	51 5 D_3_AC_LAG_3_MIN	62 9 MAX	
56	(z1:100.0D)	7 MIN	3 D_1_AC_LAG_1_MIN	51 6 D_2:3_CC_LAG_0_MIN	63 0 Property_z2_21	
57	Property_z3_10	16 D_4:3_CC_LAG_2_	3 D_1_AC_LAG_1_MIN	51 7 D_1:3_CC_LAG_9_MAX	63 1 D_3:4_CC_LAG_0_ MIN	
58	Property_z4_6	8 MAX	4 D_0_AC_LAG_3_MAX	51 8 D_1:3_CC_LAG_9_MAX	63 2 Property_z4_28	
59	Property_z2_5	9 MAX	4 D_0_AC_LAG_3_MAX	51 9 D_2:3_CC_LAG_2_ MAX	63 3 MAX	
60	Property_z4_11	17 D_0:4_CC_LAG_3_	5 D_3:0_CC_LAG_8_MIN	52 0 D_3:4_CC_LAG_9_MIN	63 4 Property_z2_23	
61	Property_z1_3	16 D_0:4_CC_LAG_3_	5 D_3:0_CC_LAG_8_MIN	52 1 D_4:2_CC_LAG_9_MIN	63 5 D_3:2_CC_LAG_9_ MAX	
62	Variable moment	17 D_0:1_CC_LAG_5_	6 Property_z1_32	52 2 D_4:2_CC_LAG_9_MIN	63 6 D_0:2_CC_LAG_5_ MAX	
63	(z4:100.0D)	17 D_4:0_CC_LAG_2_	7 D_2:4_CC_LAG_1_MAX	52 3 D_0_AC_LAG_8_MAX	63 7 D_2:3_CC_LAG_1_ MAX	
64	D_0:4_CC_LAG_2_MA	28 MIN	7 D_2:4_CC_LAG_1_MAX	52 4 D_0:4_CC_LAG_9_MIN	63 8 Property_z2_23	
65	X	17 D_3:0_CC_LAG_0_	8 D_2:1_CC_LAG_0_MIN	52 5 D_1:0_CC_LAG_4_MAX	63 9 D_3:2_CC_LAG_9_ MAX	
66	Property_z2_4	3 MAX	8 D_2:1_CC_LAG_0_MIN	52 6 D_3:4_CC_LAG_9_MIN	64 0 Property_z3_24	
67	Property_z4_7	17 D_0:4_CC_LAG_4_	9 D_1:3_CC_LAG_2_MAX	52 7 D_1:3_CC_LAG_0_MAX	64 1 D_2:3_CC_LAG_3_ MIN	
68	D_1:2_CC_LAG_0_MI	17 D_1:0_CC_LAG_2_	9 D_1:3_CC_LAG_2_MAX	52 8 D_4:2_CC_LAG_9_MIN	64 2 Property_z3_24	
69	N	7 MIN	10 D_4:1_CC_LAG_7_MIN	52 9 D_2:0_CC_LAG_9_MAX	64 3 D_2:3_CC_LAG_3_ MAX	
70	D_1:0_CC_LAG_5_MI	17 D_4_AC_LAG_7_MI	10 D_4:1_CC_LAG_7_MIN	53 0 D_3_AC_LAG_2_MAX	64 4 Property_z2_20	
71	N	8 N	11 D_3:4_CC_LAG_5_MAX	53 1 D_2:3_CC_LAG_7_MIN	64 5 Property_z5_27	
72	Property_z1_12	17 D_0_AC_LAG_2_MI	11 D_3:4_CC_LAG_5_MAX	53 2 D_2:3_CC_LAG_8_MIN		
73	Property_z1_13	9 N	12 D_1:2_CC_LAG_8_MIN	53 3 D_2:3_CC_LAG_8_MIN		
74	Property_z4_2	18 D_2:1_CC_LAG_7_	12 D_1:2_CC_LAG_8_MIN	53 4 D_2:3_CC_LAG_1_MIN		
75	D_1:4_CC_LAG_1_MA	18 D_4:1_CC_LAG_5_	13 D_4_AC_LAG_7_MAX	53 5 D_2:0_CC_LAG_9_MAX		
76	X	19 MAX	13 D_4_AC_LAG_7_MAX	53 6 D_3_AC_LAG_2_MAX		
77	Property_z3_14	18 D_3:0_CC_LAG_6_	14 D_4:1_CC_LAG_8_MAX	53 7 D_1_AC_LAG_9_MAX		
78	Property_z4_9	20 MIN	14 D_4:1_CC_LAG_8_MAX	53 8 D_3_AC_LAG_2_MAX		
79	Variable moment	18 D_1:4_CC_LAG_0_	15 Wimley-White Partitioning	53 9 D_2:3_CC_LAG_7_MIN		
80	(z5:100.0D)	18 D_3:4_CC_LAG_1_	15 (pH5.0)	53 10 D_2:3_CC_LAG_7_MIN		
81	Property_z4_1	21 MAX	16 D_4_AC_LAG_3_MIN	53 11 D_2:3_CC_LAG_8_MIN		
82	Property_z4_10	18 D_1:2_CC_LAG_6_	16 D_4_AC_LAG_3_MIN	53 12 D_2:3_CC_LAG_8_MIN		
83	N	22 MIN	17 D_1:3_CC_LAG_1_MIN	53 13 D_2:3_CC_LAG_8_MIN		
84	Property_z2_11	18 D_2:0_CC_LAG_5_	17 D_1:3_CC_LAG_1_MIN	53 14 D_2:3_CC_LAG_8_MIN		
85	Property_z3_15	23 N	18 D_2:0_CC_LAG_5_	53 15 D_2:3_CC_LAG_8_MIN		
86	Property_z4_11	24 MAX	19 D_2:1_CC_LAG_6_	53 16 D_2:3_CC_LAG_8_MIN		
87	Property_z2_12	19 D_1:2_CC_LAG_6_	20 D_2:1_CC_LAG_6_	53 17 D_2:3_CC_LAG_8_MIN		
88	N	25 MIN	20 D_2:1_CC_LAG_6_	53 18 D_2:3_CC_LAG_8_MIN		
89	Property_z3_16	20 D_2:0_CC_LAG_5_	21 D_2:1_CC_LAG_6_	53 19 D_2:3_CC_LAG_8_MIN		
90	Property_z4_12	26 MAX	21 D_2:1_CC_LAG_6_	53 20 D_2:3_CC_LAG_8_MIN		
91	Property_z2_13	21 N	22 D_2:0_CC_LAG_5_	53 21 D_2:3_CC_LAG_8_MIN		
92	Property_z3_17	22 MAX	22 D_2:0_CC_LAG_5_	53 22 D_2:3_CC_LAG_8_MIN		
93	Property_z4_13	23 N	23 D_2:1_CC_LAG_6_	53 23 D_2:3_CC_LAG_8_MIN		
94	N	24 MAX	23 D_2:1_CC_LAG_6_	53 24 D_2:3_CC_LAG_8_MIN		
95	Property_z3_18	25 N	24 D_2:0_CC_LAG_5_	53 25 D_2:3_CC_LAG_8_MIN		
96	Property_z4_14	26 MAX	24 D_2:0_CC_LAG_5_	53 26 D_2:3_CC_LAG_8_MIN		
97	Property_z2_14	27 N	25 D_2:1_CC_LAG_6_	53 27 D_2:3_CC_LAG_8_MIN		
98	Property_z3_19	28 MAX	25 D_2:1_CC_LAG_6_	53 28 D_2:3_CC_LAG_8_MIN		
99	Property_z4_15	29 N	26 D_2:0_CC_LAG_5_	53 29 D_2:3_CC_LAG_8_MIN		
100	N	30 MAX	26 D_2:0_CC_LAG_5_	53 30 D_2:3_CC_LAG_8_MIN		

71	Property_z4_10 D_4:0_CC_LAG_6_MA	18 D_4_AC_LAG_3_M 6 AX 18 D_0:3_CC_LAG_5_	30 1 D_1:4_CC_LAG_8_MAX 30 2 D_1:2_CC_LAG_2_MIN	41 6 Property_z5_31 41 7 D_4:1_CC_LAG_6_MIN	53 1 D_1_AC_LAG_8_MIN 53 2 Property_z1_30	64 D_0:2_CC_LAG_2_
72	X D_0:3_CC_LAG_2_MA	18 D_4:3_CC_LAG_8_	30 3 D_4:0_CC_LAG_9_MAX	41 8 D_4:2_CC_LAG_7_MAX	53 3 D_2:0_CC_LAG_6_MAX	64 D_2_AC_LAG_9_M
73	X	8 MAX	30 4 D_2:0_CC_LAG_4_MIN	41 9 D_3:1_CC_LAG_8_MAX	53 4 Property_z2_17	64 Property_z2_27
74	Property_z4_8	18 D_1:4_CC_LAG_5_	30 5 D_0:4_CC_LAG_4_MIN	42 0 D_1:3_CC_LAG_9_MIN	53 5 D_3:2_CC_LAG_9_MIN	64 Property_z1_27
75	Property_z1_6	9 MAX	30 6 D_0_AC_LAG_0_MAX	42 1 D_3:4_CC_LAG_9_MAX	53 6 D_1:2_CC_LAG_7_MAX	65 D_2_AC_LAG_5_MI
76	Property_z5_5 D_3:0_CC_LAG_1_MI	19 D_0:4_CC_LAG_0_	30 7 D_3_AC_LAG_1_MIN	42 2 D_1:3_CC_LAG_1_MAX	53 7 D_3_AC_LAG_9_MIN	65 D_2_AC_LAG_3_M
77	N	2 MIN	30 8 D_0:3_CC_LAG_6_MIN	42 3 D_2:1_CC_LAG_5_MAX	53 8 D_2:3_CC_LAG_7_MAX	65 D_2_AC_LAG_5_M
78	Property_z1_0	19 D_0:3_CC_LAG_3_	30 9 D_2:4_CC_LAG_5_MAX	42 4 D_0:1_CC_LAG_7_MAX	53 9 D_2:3_CC_LAG_1_MIN	65 Property_z1_26
79	D_4_AC_LAG_4_MIN	19 D_1:0_CC_LAG_7_	31 0 D_1:4_CC_LAG_6_MAX	42 5 D_0:4_CC_LAG_8_MIN	54 0 D_2:4_CC_LAG_7_MAX	65 D_2_AC_LAG_0_M
80	Property_z4_3 D_0:4_CC_LAG_2_MI	5 MIN	31 1 D_1:3_CC_LAG_4_MIN	42 6 D_1:2_CC_LAG_4_MAX	54 1 Property_z5_19	65 AX
81	N	19 D_4:1_CC_LAG_3_	31 2 D_0:3_CC_LAG_7_MAX	42 7 D_0_AC_LAG_9_MIN	54 2 D_1:0_CC_LAG_9_MAX	65 N
82	D_1:4_CC_LAG_4_MA	6 MAX	31 3 D_0:2_CC_LAG_7_MIN	42 8 D_3_AC_LAG_1_MAX	54 3 D_4:3_CC_LAG_0_MIN	65 D_2_AC_LAG_0_MI
83	X	19 D_4:0_CC_LAG_0_	31 4 D_0_AC_LAG_2_MAX	42 9 D_4:1_CC_LAG_4_MIN	54 4 D_3:2_CC_LAG_2_MIN	65 Property_z5_26
84	Property_z5_4	8 MIN	31 5 D_4_AC_LAG_9_MIN	43 0 D_1_AC_LAG_4_MAX	54 5 Property_z3_27	65 D_2_AC_LAG_4_M
85	Property_z1_34	19 D_2:0_CC_LAG_3_	31 6 D_0:1_CC_LAG_1_MAX	43 1 D_2:0_CC_LAG_9_MIN	54 6 D_2:1_CC_LAG_2_MAX	66 0 Property_z4_18
86	NetCharge@9.0	9 MIN	31 7 D_4:1_CC_LAG_3_	43 2 D_3:1_CC_LAG_7_MAX	54 7 D_2:3_CC_LAG_2_MIN	66 D_2_AC_LAG_4_MI
87	Property_z1_14	20 D_0:3_CC_LAG_0_	31 8 D_3:4_CC_LAG_6_MIN	43 3 D_1:0_CC_LAG_9_MIN	54 8 D_2:0_CC_LAG_5_MAX	66 1 N
88	N	1 MAX	31 9 D_0:1_CC_LAG_1_MAX	43 4 D_2:0_CC_LAG_9_MIN	54 9 D_2:1_CC_LAG_6_M	66 2 Property_z2_26
89	D_3:0_CC_LAG_2_MA	20 D_4:1_CC_LAG_3_	31 10 D_3:1_CC_LAG_6_MAX	43 5 D_1:0_CC_LAG_5_MAX	55 0 D_3_AC_LAG_2_MIN	66 3 N
90	X	2 MIN	31 11 D_4:0_CC_LAG_5_	43 6 D_0:1_CC_LAG_6_MAX	55 1 D_1:2_CC_LAG_9_MAX	66 4 AX
91	Property_z3_11	20 D_4:0_CC_LAG_9_M	31 12 D_1:0_CC_LAG_9_MIN	43 7 D_4:1_CC_LAG_5_MIN	55 2 Property_z3_30	66 5 Property_z1_24
92	Property_z5_3 D_1:0_CC_LAG_0_MI	3 AX	31 13 D_2:1_CC_LAG_5_MIN	43 8 D_1_AC_LAG_3_MAX	55 3 D_3:4_CC_LAG_7_MIN	66 6 D_2_AC_LAG_8_MI
93	N	20 D_0_AC_LAG_6_M	31 14 D_3_AC_LAG_6_MAX	43 9 D_4:2_CC_LAG_2_MAX	55 4 Property_z3_19	66 7 Property_z4_25
94	D_3:0_CC_LAG_2_MA	4 AX	32 0 D_3:1_CC_LAG_6_MIN	43 10 D_1:0_CC_LAG_5_MAX	55 5 D_3_AC_LAG_2_MIN	66 8 Property_z4_19
95	X	20 D_0:1_CC_LAG_6_	32 1 D_2:1_CC_LAG_6_MAX	43 11 D_0:1_CC_LAG_6_MAX	55 6 D_1:2_CC_LAG_9_MAX	66 9 D_2_AC_LAG_7_M
96	Property_z1_10	5 MIN	32 2 D_1:0_CC_LAG_3_MAX	43 12 D_4:1_CC_LAG_5_MIN	55 7 D_1_AC_LAG_3_MAX	66 AX
97	Property_z5_0 D_4:0_CC_LAG_3_MA	20 D_4:1_CC_LAG_4_	32 3 D_2:1_CC_LAG_5_MIN	43 13 D_2:3_CC_LAG_4_MIN	55 8 D_2:3_CC_LAG_4_MIN	
98	X	8 MAX	32 4 D_4:3_CC_LAG_9_MAX	43 9 Property_z3_17		
99	D_0:3_CC_LAG_0_MI	20 D_2:4_CC_LAG_3_				
100	N	9 MAX				

95	Property_z3_6	21	D_3:0_CC_LAG_3_	32	D_1:4_CC_LAG_7_MIN	44	D_3_AC_LAG_6_MIN	55	D_3_AC_LAG_9_MAX	67	Property_z5_25
	D_1:4_CC_LAG_0_MA	0	MIN	5	D_1:2_CC_LAG_5_	44	D_0_AC_LAG_7_MIN	55	D_3_AC_LAG_9_MAX	67	D_2_AC_LAG_8_M
96	X	21	D_1:2_CC_LAG_5_	32	D_0:1_CC_LAG_8_MIN	44	D_0:1_CC_LAG_7_MIN	6	Property_z1_21	1	AX
97	D_0:1_CC_LAG_3_MI	1	MIN	32	D_0:1_CC_LAG_4_MAX	44	D_3:2_CC_LAG_2_MAX	7	D_0:2_CC_LAG_3_MAX	2	N
	N	2	Property_z1_15	32	D_0:1_CC_LAG_4_MAX	44	D_0:1_CC_LAG_9_MAX	55	D_4:2_CC_LAG_7_MIN	3	Property_z2_25
98	Property_z1_7	3	Isoelectric point	8	D_3:4_CC_LAG_8_MAX	44	D_1:4_CC_LAG_9_MIN	55	D_2:1_CC_LAG_1_MAX	4	AX
	D_0:4_CC_LAG_6_MA	21	D_3:4_CC_LAG_4_	32	D_2:0_CC_LAG_8_MAX	44	D_1:4_CC_LAG_9_MIN	9	D_2:1_CC_LAG_1_MAX	67	D_2_AC_LAG_1_M
99	X	4	MAX	33	D_2:0_CC_LAG_5_MIN	44	D_3:1_CC_LAG_4_MAX	0	D_3:0_CC_LAG_9_MAX	5	Property_z4_24
10	D_0:4_CC_LAG_0_MI	21	D_0:1_CC_LAG_2_	33	D_2:0_CC_LAG_5_MIN	44	D_3:2_CC_LAG_1_MIN	1	D_4:3_CC_LAG_9_MIN	67	D_2_AC_LAG_2_M
0	N	5	MIN	33	D_0:3_CC_LAG_6_MAX	44	D_2:0_CC_LAG_4_MAX	2	D_4:2_CC_LAG_3_MIN	7	Property_z4_20
10	Property_z3_15	21	D_4:0_CC_LAG_8_	33	D_0:3_CC_LAG_6_MAX	44	D_1:0_CC_LAG_8_MAX	3	D_3:2_CC_LAG_1_MAX	8	Property_z5_24
10	D_1:2_CC_LAG_4_MI	6	MAX	33	D_1_AC_LAG_7_M	44	D_0:3_CC_LAG_31	4	Property_z5_17	9	N
2	N	7	AX	33	D_1_AC_LAG_0_MIN	45	D_3_AC_LAG_5_MIN	5	D_1:2_CC_LAG_3_MAX	0	N
10	D_0:4_CC_LAG_5_MA	21	D_1:3_CC_LAG_2_	33	D_3:1_CC_LAG_5_MIN	45	D_2:1_CC_LAG_8_MIN	6	Property_z5_29	1	Property_z5_23
3	X	8	MIN	33	D_3:1_CC_LAG_5_MIN	45	D_2:4_CC_LAG_8_MAX	7	D_2:3_CC_LAG_6_MAX	2	Property_z4_23
10	Property_z4_13	21	D_1:2_CC_LAG_1_	33	D_1:2_CC_LAG_6_MAX	45	D_3:1_CC_LAG_3_MAX	8	D_4:2_CC_LAG_4_MIN	3	Property_z1_25
5	Property_z3_4	9	Property_z3_33	4	D_1:2_CC_LAG_6_MAX	45	Wimley-White Partitioning	4	D_2:4_CC_LAG_6_MIN	4	Property_z4_21
10	D_4:1_CC_LAG_2_MA	22	D_1:0_CC_LAG_1_	33	D_1_AC_LAG_2_MAX	45	(pH7.0)	5	D_0:2_CC_LAG_4_MAX	5	Property_z4_22
6	X	1	MAX	6	Property_z2_15	45	D_3_AC_LAG_0_MAX	0	Property_z5_20	6	Molecular Weight
10	Property_z2_33	22	D_4:3_CC_LAG_6_	33	Property_z2_15	45	D_4:1_CC_LAG_9_MIN	1	D_2:3_CC_LAG_5_MAX	7	Property_z5_22
7	D_0:3_CC_LAG_1_MI	2	Property_z4_14	7	D_0_AC_LAG_9_MAX	45	D_4:2_CC_LAG_1_MAX	2	D_2:4_CC_LAG_9_MIN	8	Size
10	N	22	D_1:4_CC_LAG_3_	33	D_3_AC_LAG_5_MAX	45	D_4:2_CC_LAG_1_MAX	3	D_2:4_CC_LAG_9_MIN	9	Property_z5_21
8	D_4:1_CC_LAG_0_MI	3	MAX	8	D_0_AC_LAG_5_M	45	D_1:3_CC_LAG_6_MAX	4	D_0:2_CC_LAG_9_MAX	5	Property_z4_29
9	N	4	AX	9	D_1:4_CC_LAG_2_MIN	45	D_3:0_CC_LAG_9_MIN	5	D_3:2_CC_LAG_3_MIN	6	D_4:2_CC_LAG_2_MIN
11	Property_z2_33	22	D_3:0_CC_LAG_7_	34	Property_z2_32	45	D_4:1_CC_LAG_9_MIN	6	D_3:2_CC_LAG_3_MIN	7	Property_z4_29
0	D_4_AC_LAG_6_MIN	5	MIN	0	D_4_AC_LAG_5_M	45	D_4:2_CC_LAG_1_MAX	7	D_4:2_CC_LAG_7_MIN	8	Property_z1_25
11	D_0:2_CC_LAG_1_MI	22	D_4_AC_LAG_5_M	34	D_1:2_CC_LAG_1_MIN	45	D_1:3_CC_LAG_6_MAX	8	D_2:4_CC_LAG_9_MIN	9	Size
1	N	6	AX	1	D_4:2_CC_LAG_4_MAX	45	D_3:0_CC_LAG_9_MIN	9	D_0:2_CC_LAG_9_MAX	6	Property_z5_21
11	D_4:3_CC_LAG_5_MA	22	D_4:3_CC_LAG_1_	34	D_4:2_CC_LAG_4_MAX	45	D_4:2_CC_LAG_1_MAX	1	D_3:2_CC_LAG_3_MIN	5	Property_z4_29
2	X	7	MAX	2	D_2:1_CC_LAG_0_	45	D_1:3_CC_LAG_6_MAX	2	D_2:3_CC_LAG_5_MAX	6	D_4:2_CC_LAG_2_MIN
11	D_1:0_CC_LAG_6_MI	22	D_2:1_CC_LAG_0_	34	D_3:0_CC_LAG_9_MIN	45	D_2:4_CC_LAG_9_MIN	3	D_2:4_CC_LAG_9_MIN	7	Property_z5_22
3	N	8	MAX	3	D_3:1_CC_LAG_0_	45	D_3:0_CC_LAG_4_MAX	4	D_0:2_CC_LAG_9_MAX	8	Size
11	D_0_AC_LAG_0_MIN	22	D_3:1_CC_LAG_0_	34	D_3:0_CC_LAG_4_MAX	45	D_3:0_CC_LAG_4_MAX	5	D_3:2_CC_LAG_3_MIN	9	Property_z5_21
11	Property_z3_5	9	MIN	4	D_4:1_CC_LAG_7_	46	Property_z3_29	6	D_4:2_CC_LAG_2_MIN	6	D_4:2_CC_LAG_2_MIN
11	D_3:4_CC_LAG_3_MA	0	MAX	5	D_0:3_CC_LAG_7_MIN	46	D_0:3_CC_LAG_7_MIN	5	D_3:2_CC_LAG_3_MIN	7	Property_z5_22
6	X	23	D_3:4_CC_LAG_6_	34	D_0:3_CC_LAG_7_MIN	46	D_0:3_CC_LAG_5_MAX	6	D_4:2_CC_LAG_2_MIN	8	Size
		1	MAX	6	D_0:3_CC_LAG_5_MAX	46	Property_z4_15	1	D_4:2_CC_LAG_2_MIN	9	Property_z5_21

2.2 Figures

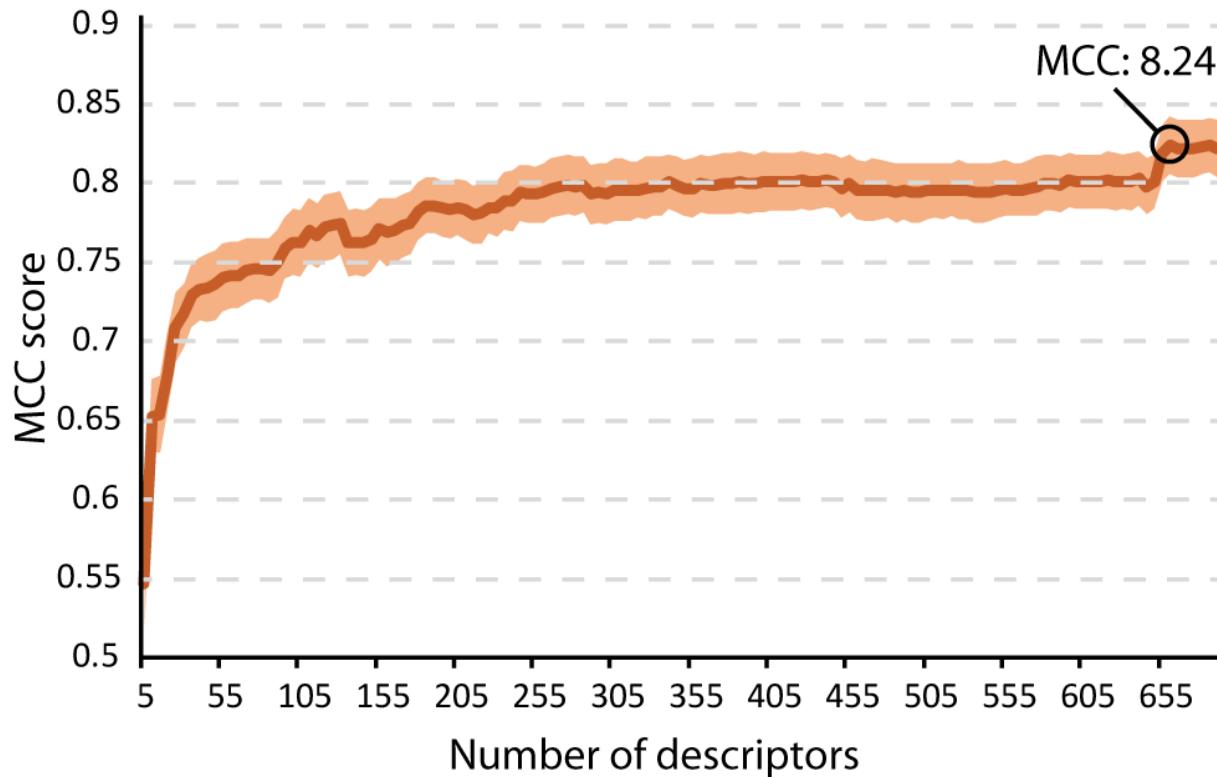


Figure S1. IFS results. Ten-fold cross validation of the sorted list of descriptors. The descriptor list was sorted by Maximum Relevance, Minimum Redundancy (mRMR) and a total number of 138 models were trained. The model giving the highest MCC score was selected.