

## Supplementary data:

### **Bitter or not? BitterPredict, a tool for predicting taste from chemical structure**

Ayana Dagan-Wiener<sup>1,2</sup>, Ido Nissim<sup>1,2</sup>, Natalie Ben Abu<sup>1,2</sup>, Gigliola Borgonovo<sup>3</sup>, Angela Bassoli<sup>3</sup>, Masha Y. Niv<sup>1,2\*</sup>

<sup>1</sup>Institute of Biochemistry, Food Science and Nutrition, The Robert H. Smith Faculty of Agriculture, Food, and Environment, The Hebrew University of Jerusalem, Rehovot, 76100, Israel.

<sup>2</sup>The Fritz Haber Center for Molecular Dynamics, The Hebrew University of Jerusalem, 91904, Israel.

<sup>3</sup> DeFENS-Department for Food, Environmental and Nutritional Sciences, University of Milan, Via Celoria 2, Milano, 20133, Italy.

\*correspondence to [masha.niv@mail.huji.ac.il](mailto:masha.niv@mail.huji.ac.il)

### **Table S1 Additional Supplementary files:**

File Name	Description
validation.xls	1) Holds structures (smiles format), Bitter label (0-non-bitter , 1 bitter) and BitterPredict results for the three external sets used for validation (Bitter New, UNIMI set, Phyto. Dictionary). 2) Holds the Literature mining validation on the top 30 predicted bitter and non-bitter from DrugBank. 3) Holds the predicted bitter compounds (score>0.6) and predicted non-bitter compounds (score<-0.7) from Sigma-Aldrich Ingredients Catalog Flavors&Fragrances food with explanation on top predicted compounds that were not used for validation .
prospective_prediction_sets.xls	Holds structures (smiles format) and BitterPredict results for the four datasets used for the prospective predictions (FooDB, DrugBank,ChEBI, Natural Products ZINC15)

### **Table S2: Prediction results using only the non-bitter flavors as negative group.**

Results of AdaBoost model which was trained only with non-bitter flavors as negative set.

Table 1.A shows good performance measure on the hold-out test set.

Table 1.B shows the high percentage predicted to be bitter in FooDB DrugBank and ChEBI.

A.	Sensitivity TP/TP+FN	Specificity TN/TN+FP	Accuracy (TP+TN)/Total
AdaBoost	0.75	0.88	0.84
B.	FooDB	DrugBank	ChEBI (random)
% molecules from set predicted to be Bitter	43.2%	80.8%	70.3%

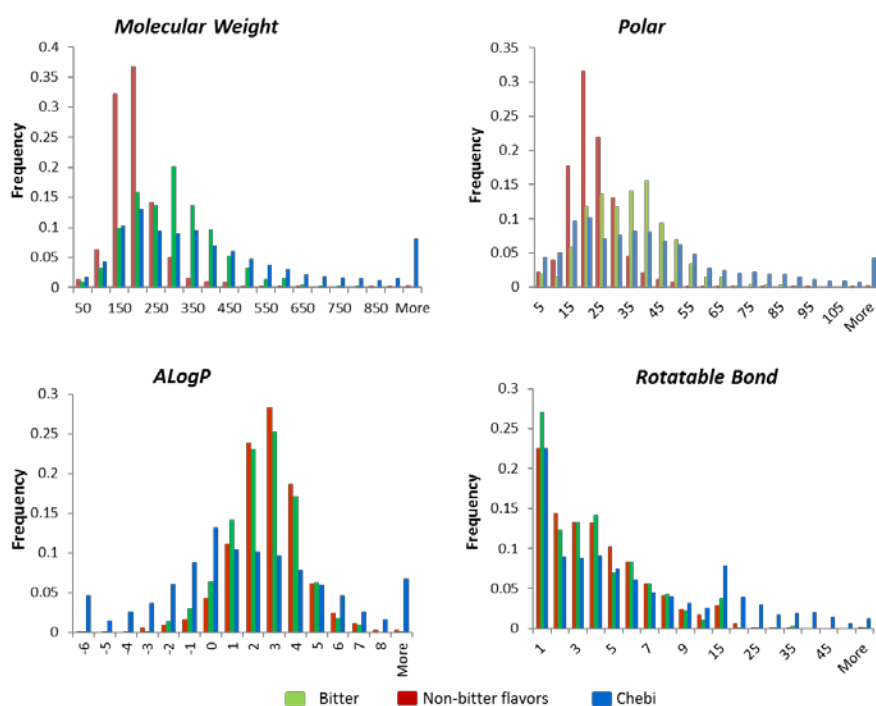
TP- number of positive (bitter) molecules correctly classified.

FP- number of negative (non-bitter) molecules incorrectly classified.

TN- number of negative (non-bitter) molecules correctly classified.

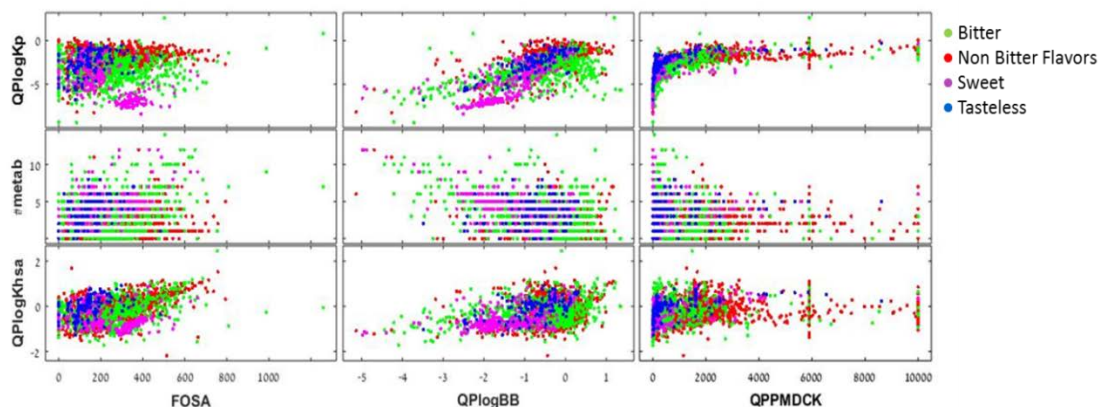
FN- number of positive molecules (bitter) incorrectly classified.

**Figure S1:**



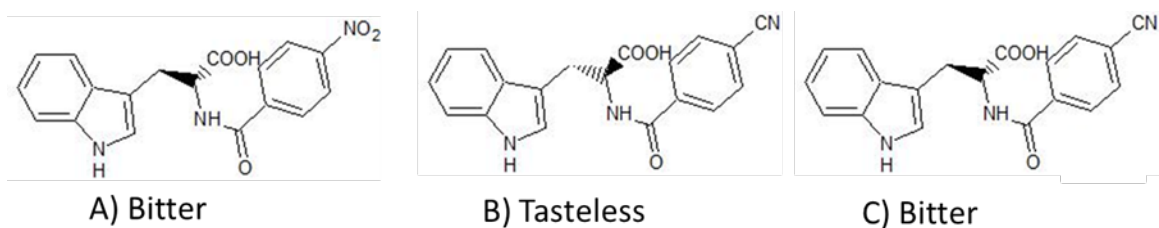
**Figure S1: Comparison of selected physicochemical properties between the Bitter, Non-bitter flavors and ChEBI sets.** In terms of The Molecular weight and Polarity the Non-bitter flavors set has low variance and occupies narrower range compared to the Bitter and ChEBI set. The hydrophobicity of the Bitter and Non-bitter flavors set, represented by AlogP, is similar and somewhat higher compared to random molecules (ChEBI). The distribution of number of rotatable bonds is similar between the three sets.

**Figure S2:**

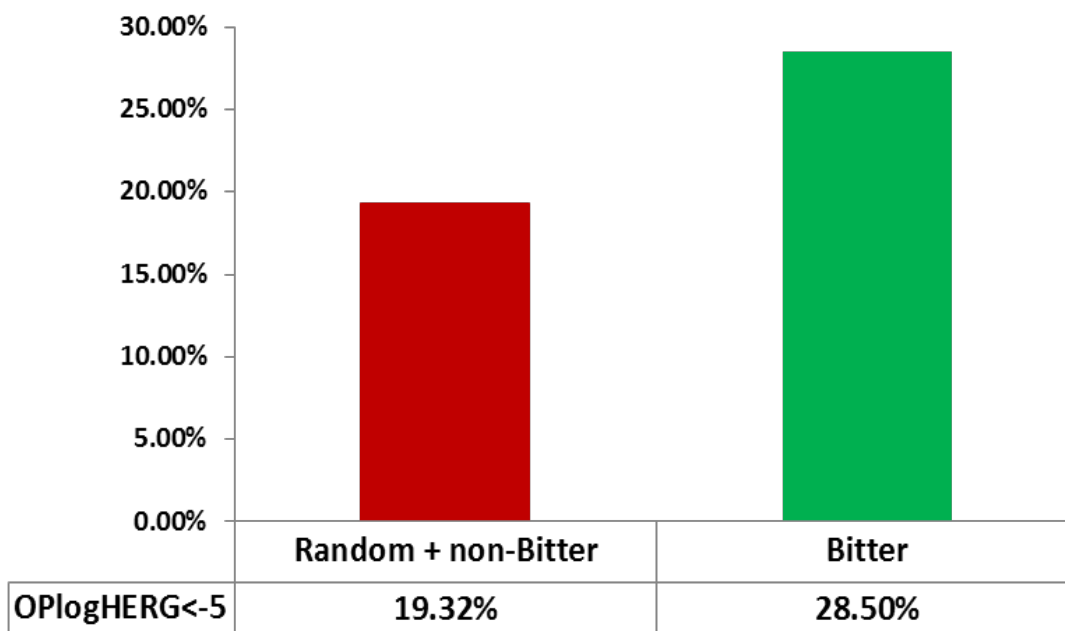


**Figure S2: Bivariate statistical analysis of selected QikProp properties.** The bivariate statistical analysis indicates combined properties ranges which are enriched with bitter molecules. Most of the bitter molecules properties tend to reside in the range that is suggested by QikProp (version 4.6 Schrödinger, LLC, New York, NY, 2015) as acceptable for drugs  $-3 < \text{QPlogBB} < 1.2$ ,  $-1.5 < \text{Prediction of binding to human serum albumin (QPlogKhsa)} < 1.5$ ,  $-8 < \text{Predicted skin permeability (QPlogKp)} < -1$ , Predicted apparent Madin-Darby Canine Kidney Epithelial Cells (MDCK) cell permeability in nm/sec.  $(\text{QPPMDCK}) > 25$ , Number of likely metabolic reactions  $0 < (\# \text{metab}) < 8$ , Hydrophobic component of the total solvent accessible surface area  $(\text{FOSA}) < 750$  (for more details about the suggested ranges please refer to QikProp manual: [http://gohom.win/ManualHom/Schrodinger/Schrodinger\\_2015-2\\_docs/qikprop/qikprop\\_user\\_manual.pdf](http://gohom.win/ManualHom/Schrodinger/Schrodinger_2015-2_docs/qikprop/qikprop_user_manual.pdf))

**Figure S3:**



**Figure S3: UNIMI set molecules:** Examples of 3 molecules from the challenging UNIMI set, demonstrating that similar, molecules (A,B) and even stereo isomers (B,C) elicit different tastes. BittePredict predicts all three molecules as bitter, though B and C got relatively low bitterness score (<0.2)



**Figure S4: Predicted IC50 value for blockage of Human ether-a-go-go-related gene (hERG) potassium channels in the bitter and random+non bitter sets.** We found that bitter compounds are enriched (~30%) with compounds which are predicted to block gene (hERG) K+ channels compared to random + non-bitter set (20%). The result fits with our recent finding that >20% of the bitter tastants in BitterDB are reported to inhibit hERG. It was recently found that bitter receptors are also expressed in heart tissue, which might be connected to this common off- target of bitter molecules.

**Table S3 descriptors used in BitterPredict classifier:**

<p>47 descriptors from QikProp</p>	<p><b>#stars</b>- Number of property or descriptor values out of 24 properties and QikProp descriptors that fall outside the 95% range of similar values for known drugs.</p> <p><b>#amine</b> Number of non-conjugated amine group</p> <p><b>#amidine</b> Number of amidine and guanidine groups.</p> <p><b>#acid</b> Number of carboxylic acid groups.</p> <p><b>#amide</b> Number of non-conjugated amide groups.</p> <p><b>#rotor</b> Number of non-trivial (not CX3), non-hindered (not alkene, amide, small ring) rotatable bonds.</p> <p><b>#rtvFG</b> Number of reactive functional groups;</p> <p><b>CNS</b> - Predicted central nervous system activity</p> <p><b>dipole†</b> Computed dipole moment of the molecule.</p> <p><b>SASA</b> Total solvent accessible surface area (SASA) in square angstroms using a probe with a 1.4 Å radius.</p> <p><b>FOSA</b> Hydrophobic component of the SASA (saturated carbon and attached hydrogen).</p> <p><b>FISA</b> Hydrophilic component of the SASA (SASA on N, O, H on heteroatoms, carbonyl C). 7.0 – 330.0</p> <p><b>PISA</b> π (carbon and attached hydrogen) component of the SASA.</p> <p><b>WPSA</b> Weakly polar component of the SASA (halogens, P, and S).</p> <p><b>volume</b> Total solvent-accessible volume in cubic angstroms using a probe with a 1.4 Å radius.</p> <p><b>donorHB</b> Estimated number of hydrogen bonds that would be donated by the solute to water molecules in an aqueous solution.</p> <p><b>accptHB</b> Estimated number of hydrogen bonds that would be accepted by the solute from water molecules in an aqueous solution</p> <p><b>dip<sup>2</sup>/V†</b> Square of the dipole moment divided by the molecular volume.</p> <p><b>ACxDN<sup>.5</sup>/SA</b> Index of cohesive interaction in solids.</p> <p><b>glob</b> Globularity descriptor.</p> <p><b>QPpolrz</b> Predicted polarizability in cubic angstroms.</p> <p><b>QPlogPC16</b> Predicted hexadecane/gas partition coefficient.</p> <p><b>QPlogPoct</b> Predicted octanol/gas partition coefficient.</p> <p><b>QPlogPw</b> Predicted water/gas partition coefficient.</p> <p><b>QPlogPo/w</b> Predicted octanol/water partition coefficient.</p> <p><b>QPlogS</b> Predicted aqueous solubility</p> <p><b>CIQPlogS</b> Conformation-independent predicted aqueous solubility,</p> <p><b>QPlogHERG</b> Predicted IC50 value for blockage of HERG K+ channels.</p>
------------------------------------	---

	<p><b>QPPCaco</b> Predicted apparent Caco-2 cell (Caco- 2 cells are a model for the gut-blood barrier, (non-active transport).</p> <p><b>QPlogBB</b> Predicted brain/blood partition coefficient.</p> <p><b>QPPMDCK</b> Predicted apparent MDCK (Madin-Darby Canine Kidney Epithelial) cell permeability in nm/sec.</p> <p><b>QPlogKp</b> Predicted skin permeability.</p> <p><b>IP(ev)</b> calculated ionization</p> <p><b>EA(eV)</b> calculated electron affinity</p> <p><b>#metab</b> Number of likely metabolic reactions.</p> <p><b>QPlogKhsa</b> Prediction of binding to human serum albumin.</p> <p><b>HumanOralAbsorption</b> Predicted qualitative human oral absorption</p> <p><b>PercentHumanOralAbsorption</b> Predicted human oral absorption on 0 to 100% scale.</p> <p><b>SAFluorine</b> Solvent-accessible surface area of fluorine atoms</p> <p><b>SAamideO</b> Solvent-accessible surface area of amide oxygen atoms.</p> <p><b>PSA</b> Van der Waals surface area of polar nitrogen and oxygen atoms and carbonyl carbon atoms.</p> <p><b>#NandO</b> Number of nitrogen and oxygen atoms.</p> <p><b>RuleOfFive</b> Number of violations of Lipinski's rule of five. The rules are: mol_MW &lt; 500, QPlogPo/w &lt; 5, donorHB ≤ 5, accptHB ≤ 10.</p> <p><b>RuleOfThree</b> Number of violations of Jorgensen's rule of three. The three rules are: QPlogS &gt; -5.7, QP PCaco &gt; 22 nm/s, # Primary Metabolites &lt; 7.</p> <p><b>#ringatoms</b> Number of atoms in a ring</p> <p><b>#in34</b> Number of atoms in 3- or 4-membered rings</p> <p><b>#in56</b> Number of atoms in 5- or 6-membered rings</p> <p><b>#noncon</b> number of ring atoms not able to form conjugated aromatic systems (e.g. sp<sup>3</sup> C).</p> <p><b>#nonHatm</b> Number of heavy atoms (non hydrogen atoms)</p>
<p>12 <i>Physiochemical</i> and topological descriptors calculated by Canvas and QikProp)</p>	<p><b>MW</b> molecular weight,  <b>AlogP</b> lipophilicity,  <b>RB</b> rotatable bonds count,  <b>PSA</b> polar surface area,  <b>Estate</b> electrotopological states,  <b>MR</b> molecular refractivity,  <b>Polar</b> molecular polarizability ,  <b>Ar ring</b> aromatic rings count ,  <b>Ring</b> rings count,  <b>Chiral</b> chiral centers count ,  <b>HA</b> heavy atoms count ,  <b>Total charge</b>.</p>

