

THE PREDICTION OF pK_a VALUES FOR PHENOLIC COMPOUNDS BY THE DFT THEORY

Nguyen Thi My, Nguyen Van Din, Mai Van Bay*

The University of Danang - University of Science and Education

*Corresponding author: mvbay@ued.udn.vn

(Received: January 25, 2022; Accepted: March 24, 2022)

Abstract - The acid dissociation constant is an important parameter that affects the physicochemical properties of molecules in solution. A set of 20 phenolic compounds were used to establish a model to predict pK_a values of phenolic compounds. Calculations in aqueous medium were performed with a polarizable continuum solvent model (PCM) and three hybrid DFT functionals (B3LYP, PBE0, ω B97XD and M062X) with the basis set 6-311++G(d,p). The directly calculated value of pK_a gives less accurate results with an average absolute error (MAE) of 1.74 pK_a units when using the ω B97XD functional and phenol as reference compound. In the case of using statistical correction, the accuracy of pK_a is greatly improved. In the case of using statistical correction, the accuracy of pK_a is greatly improved with the lowest MAE value of 0.14 pK_a units (M062X; $R^2 = 0.978$). The calculated results of pK_a in this study have the same accuracy as the experimental measurements.

Key words - pK_a ; density functional theory methods (DFT); phenolic.

1. Introduction

The proton transfer reactions are one of the most basic and common types of reactions in chemistry [1]. The ability of a substance to transfer protons in a medium is characterized by its acid dissociation constant (K_a) and is usually reported as a pK_a value. The many chemical compounds act as Brønsted–Lowry acids or bases in aqueous media. Therefore, depending on the pK_a and pH values, these compounds can be ionized to different degrees and thus determine their forms of existence in the aquatic environment. The bioactive molecules from nature as well as drug molecules are usually weak acids or weak bases, so their degree of ionization in the medium affects lipophilicity, solubility, protein binding and mobility across the plasma membrane and thus pK_a affects the absorption, distribution, metabolism, excretion and toxicity properties of compounds [2]. Therefore, the pK_a value is a very important parameter for pharmaceutical compounds as well as other commonly used compounds. Most of the pK_a values are determined experimentally. However, the number of experimental pK_a values is too small compared to the extremely abundant number of chemical compounds. Therefore, there have been many publications using different theoretical approaches to predict pK_a value such as using the QSAR model (quantitative structure – activity relationship) combined with machine learning method [3, 4]; The semiempirical quantum mechanical methods and the density functional theory methods (DFT) [5-7]. The results of these publications show that it is possible to predict pK_a from the theory. However, each family of compounds needs to be approached separately, and the

quantitative results of each model must be based on compounds with similar chemical structures.

The subjects in this study were phenolic compounds. They are compounds with one or more hydroxyl groups attached to the aromatic benzene nucleus and are found in most plant tissues [8]. Phenolics are known to have biological activities, such as antibacterial and antiviral properties, anti-inflammatory and antiproliferative activities and especially many phenolics have strong antioxidant activity and are the most abundant source of antioxidants in the human diet [9].

In aqueous media, phenolics are weak acids, which is due to the proton dissociation of the OH group attached to the benzene ring. Depending on the pK_a value and pH of the body's environment, phenolics can exist as neutral molecules or ionic form, or both. Therefore, the acid-base equilibrium of phenolics in the environment affects their bioactivity [10, 11] and pK_a values of phenolics are of the utmost importance in practical applications. In this study, we use DFT theory combined with a linear regression model to predict pK_a of phenolic compounds based on a sample of 20 phenolic compounds with known pK_a from the experiment (Figure 1).

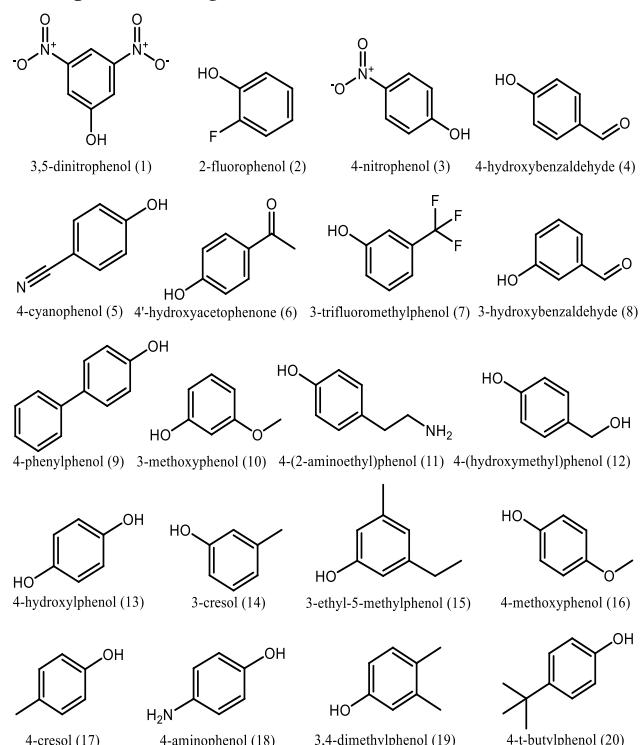


Figure 1. Phenolic compounds used in this study

2. Theoretical and computational methods

2.1. Calculation of pKa

The dissociation of phenolics (ArOH) in aqueous medium (aq) is expressed in terms of Equation (1).



Acid dissociation constant (K_a) and $\text{p}K_a$ are defined according to equations (2) and (3) respectively[12].

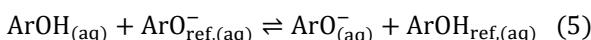
$$K_a = [\text{H}_3\text{O}^+][\text{ArO}^-]/[\text{ArOH}] \quad (2)$$

$$\text{p}K_a = -\log K_a \quad (3)$$

The value of K_a is calculated based on Equation (4).

$$\Delta G_T^0(1) = -RT \ln K_a \quad (4)$$

Where, $\Delta G_T^0(1)$ is the change in Gibbs energy of acid dissociation, calculated by DFT theory. The disadvantage of this method is that it is difficult to accurately calculate the Gibbs energy of the H_3O^+ ion because the ion is very strongly hydrated in solution[13]. This can be overcome by using an acid reference (ArOH_{ref}). Consider the following equilibria:



The $\text{p}K_a$ value of ArOH is determined according to Equation (6) (detailed description in section 1.1 of supplementary information).

$$\text{p}K_a = \frac{\Delta G_T^0(5)}{RT \ln 10} + \text{p}K_a^{\text{ref}} \quad (6)$$

Where, $\text{p}K_a^{\text{ref}}$ is the experimental $\text{p}K_a$ value the reference phenolics, ArOH_{ref} and $\Delta G_T^0(5)$ is the change in Gibbs energy of reaction on Equation (5), determined by:

$$\Delta G_T^0(5) = G_T^0(\text{ArOH}_{\text{ref}}) + G_T^0(\text{ArO}^-) - G_T^0(\text{ArOH}) - G_T^0(\text{ArO}_{\text{ref}}^-) \quad (7)$$

The study used 3 reference phenolic compounds with significantly different $\text{p}K_a$ values including phenol ($\text{p}K_a = 9.99$), 2-hydroxybenzaldehyde ($\text{p}K_a = 8.37$) and 2-cyanophenol ($\text{p}K_a = 6.86$). The predicted value of $\text{p}K_a$ by direct calculation method (denoted by $\text{p}K_a^{\text{calc}}$) is calculated according to Equation (3).

The results are statistically processed as follows: from the $\text{p}K_a^{\text{calc}}$ values, build a linear regression equation according to Equation (8).

$$\text{p}K_a^{\text{exp}} = a + b \cdot \text{p}K_a^{\text{calc}} \quad (8)$$

Where, $\text{p}K_a^{\text{exp}}$ is the experimental $\text{p}K_a$ of phenolics used in the study sample. The final predicted $\text{p}K_a$ after statistical correction (denoted by $\text{p}K_a^{\text{corr}}$) is calculated in terms of Equation (9).

$$\text{p}K_a^{\text{corr}} = \alpha + \beta \cdot \text{p}K_a^{\text{calc}} \quad (9)$$

Where, α và β are equal to the values of a and b in Equation (8), respectively.

2.2. The DFT calculation

In this study, four density functional theory (DFT) methods including B3LYP [14], PBE0 [15], ω B97XD [16] and M06-2X [17] combined with the 6-311++G(d,p) basic set were recommended for predict $\text{p}K_a$ values. The choice of these methods was based on the previous studies

because of its accuracy and lower computational cost [18-22]. In addition, the solvent effects were modeled by the polarizable continuum model (PCM) which is the most commonly used solvation model [23, 24].

The stable geometry of the compounds was checked by harmonic frequency calculation to confirm the structures which were minimal on the PES surface and to obtain the thermal and entropy contributions to the Gibbs energy. All calculations were performed at standard conditions in the solution of 1 M and 298.15 K and used Gaussian 16 Revision A.03 software[25].

3. Results and discussion

3.1. Direct calculation of $\text{p}K_a$

The $\text{p}K_a$ error between calculation and experiment of 20 phenol compounds when using different DFT functionals and reference compounds is presented in Table S2 (supplementary information: SI) and visualized in Figure 2. The results show that the error $\text{p}K_a$ is less dependent on the type of used functional but strongly depends on the reference compound. The 20 phenolic compounds in Figure 2 are numbered in order of increasing $\text{p}K_a$ value (see Table S1 in SI). The errors tend to increase as the $\text{p}K_a$ of the reference compound is further away from the $\text{p}K_a$ of the calculated compound. Specifically, the reference compound 2-cyanophenol with a small $\text{p}K_a$ (6.86) will cause large errors for the compound with a large $\text{p}K_a$ (compound 7 to 20 have a $\text{p}K_a$ between 8.98 and 10.39). In contrast, phenol with large $\text{p}K_a$ (9.99) produces small errors for compounds with similar $\text{p}K_a$ but large errors for compounds with small $\text{p}K_a$ (compound 1 to 6 with $\text{p}K_a$ of 6.69 to 8.05). The case of 2-hydroxybenzaldehyde has a $\text{p}K_a$ (8.37) lower than phenol and greater than 2-cyanophenol, so it generally does not cause too large errors in the regions with large and small $\text{p}K_a$ in the sample of the studied compounds. A general trend that does not depend on the functionals or the used reference compounds is that the error is mostly negative for substances with small $\text{p}K_a$ and positive for compounds with large $\text{p}K_a$.

The results in Figure 3 show that the mean absolute error is the smallest when using phenol as the reference compound and the largest when using 2-cyanophenol. This is explained by the fact that the majority of phenolics have a $\text{p}K_a$ closer to phenol than 2-cyanophenol. For the sample of compounds studied, the best calculation $\text{p}K_a$ results were observed by the combination of ω B97XD functional and phenol as the reference compound with a MAE value of 1.74. However, the error between calculated and experimental $\text{p}K_a$ varied widely from -5.86 (compound 1) to 2.02 (compound 18). The results here indicate that the selection of a reference compound is very important in the theoretical prediction of $\text{p}K_a$ value. The closer the $\text{p}K_a$ of the reference compound is to the one of interests, the more accurate the result. In general, the $\text{p}K_a$ calculated by the direct method using the reference compound has a large error compared with the experimental method, and therefore caution should be exercised when using this method to predict the $\text{p}K_a$ of phenolic compounds.

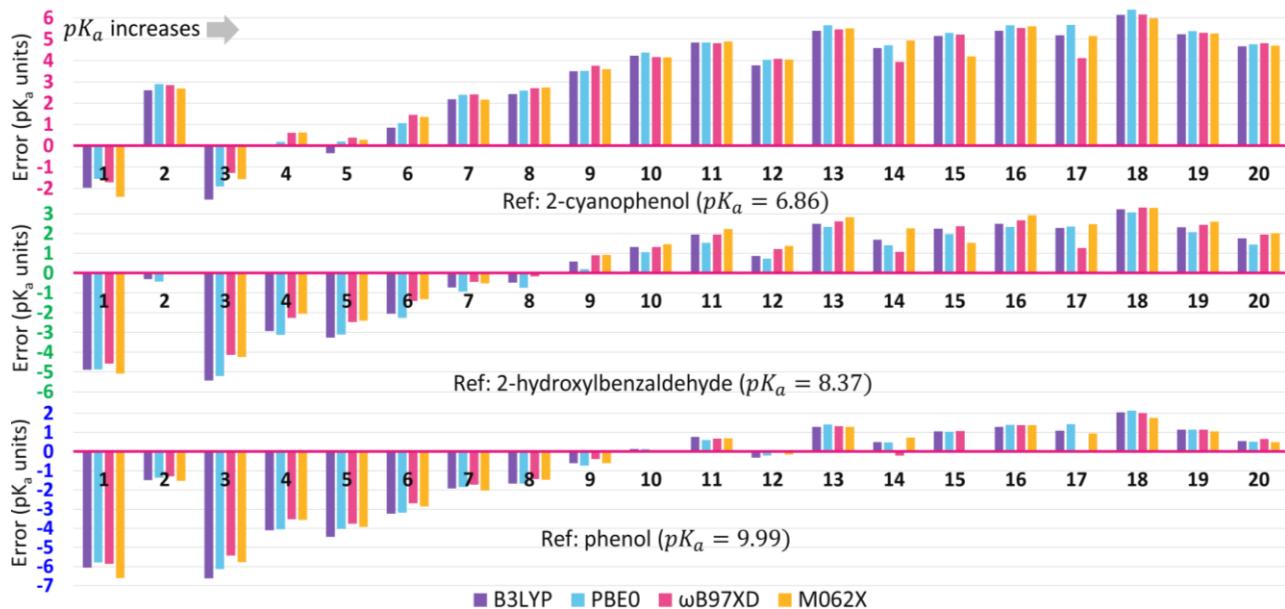


Figure 2. The error differences between pK_a calculated and experimental values without statistical correction

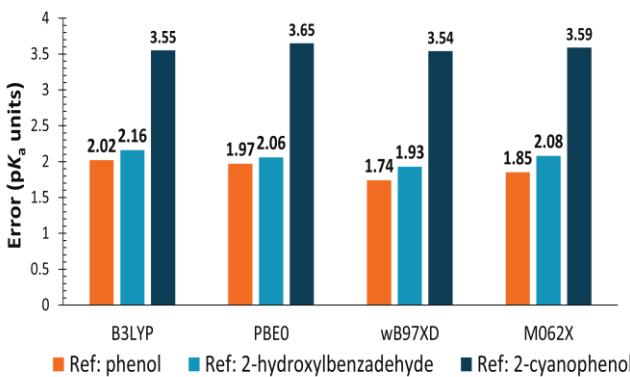


Figure 3. The mean absolute error (MAE) differences between pK_a calculated and experimental values without statistical correction

3.2. The Calculated pK_a by statistical correction

The results in section 3.1 show that the pK_a when calculated by the direct method has a large uncertainty and is highly dependent on the reference compound. In this part of the study, we will examine the correlation between the calculated and experimental pK_a values. The results presented in Figure 4 show that there is a strong linear correlation between the calculated (pK_a^{calc}) and experimental pK_a (pK_a^{exp}) with the R^2 coefficient in the range from 0.972 to 0.978. Furthermore, the correlation coefficients and slopes in the regression equations do not depend on the type of reference compound selected. Therefore, the choice of the reference compound does not affect the degree of correlation between the calculated and experimental pK_a values. This means that the error in calculating pK_a after statistical correction will not depend on the type of reference compound and this will be analyzed below.

Based on the correlation between the calculated and experimental pK_a , the predicted pK_a value is statistically

corrected (denoted by pK_a^{corr}) by replacing the quantity pK_a^{exp} in the regression equations with pK_a^{corr} . The equations for calculating pK_a^{corr} when using phenol reference compounds are presented in Table 1. Here, it should be noted that the choice of the reference compound is just a procedure when calculating, has absolutely no effect on the value of the pK_a^{corr} and is therefore the same when using the regression equations of the other reference compounds for the investigation.

Table 1. Equations for calculation of pK_a^{corr} according to pK_a^{calc} of DFT functionals when phenol is reference

Functionals	Equations	R^2
B3LYP	$pK_a^{corr} = 0.306 \times pK_a^{calc} + 6.715$	0.978
PBE0	$pK_a^{corr} = 0.314 \times pK_a^{calc} + 6.450$	0.978
wB97XD	$pK_a^{corr} = 0.333 \times pK_a^{calc} + 6.617$	0.972
M062X	$pK_a^{corr} = 0.321 \times pK_a^{calc} + 6.099$	0.978

The error between pK_a^{corr} and pK_a^{exp} of 20 phenol compounds when using different DFT functionals is presented in Table S3 (SI) and visualized in Figure 5. The results show that the errors do not depend on the reference compound. Thus, unlike the case of direct pK_a calculation, the error here tends to be uniformly distributed throughout the sample where the pK_a varies over a wide range (from as low as 6.69 for compound 1 to a maximum of 10.39 for compound 20). The most important thing is that the error in calculating pK_a with statistical correction gives much better results than in the case of calculating pK_a directly. The maximum error of the studied phenolics was only 0.4 pK_a units for compound 17 when using the wB97XD functional, while the largest error without statistical correction was 6.61 for compound 1 when using the M062X functional combined with phenol as the reference compound. As discussed, the error of the pK_a value after statistical correction is independent of the type of the used reference compound.

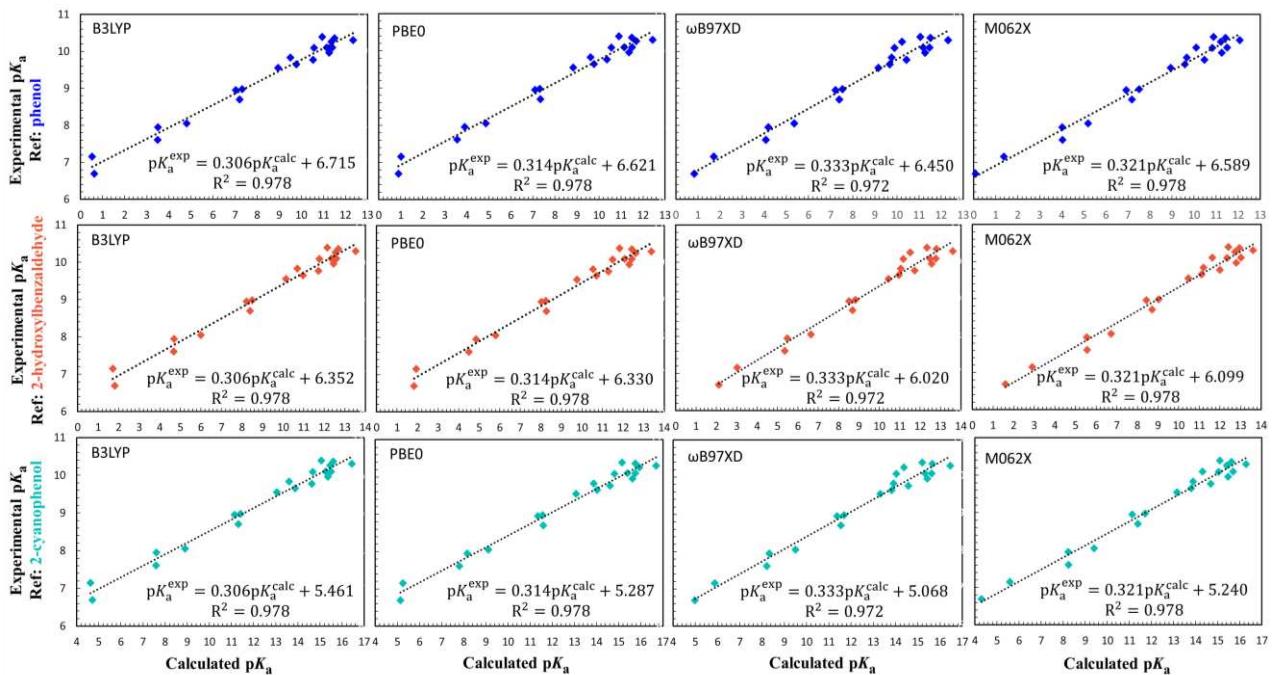


Figure 4. Diagrams of linear correlation between the pK_a^{calc} and pK_a^{exp}



Figure 5. The error differences between pK_a calculated and experimental values with statistical correction

The results in Figure 6 show that there is a very small difference in the MAE values when using the B3LYP, PBE0, ω B97XD and M062X functionals. In which, the M062X functional gives the smallest MAE value of 0.14 p K_a units and the largest absolute error when using the M062X functional observed in compound 4 is 0.27 p K_a units. In summary, by using appropriate functionals combined with statistical correction, we are able to calculate the p K_a values of phenolic compounds with an accuracy asymptotic to the experimental measurements.

4. Conclusions

This study presents a reliable procedure for calculating the p K_a of phenolic compounds in aqueous medium. The study sample consisted of 20 phenolic compounds and 3 reference phenolic compounds with experimentally known p K_a . Three hybrid DFT functionals (B3LYP, PBE0, ω B97XD and M062X) with the basis set 6-311++G(d,p) and the polarizable continuum solvent model (PCM) were used to calculate p K_a . The p K_a values calculated directly without statistical processing has a large error ($MAE \geq 1.74$ p K_a units) and the accuracy of the calculation depends strongly on the reference compound. The more accurate the result, the closer the reference compound has a p K_a to the compound of interest. In the sample of the compound studied, the ω B97XD functional combined with phenol as the reference compound gave the best results ($MAE=1.74$). ($MAE = 1.74$). The p K_a calculation results are greatly improved after statistical processing ($MAE = 0.14$ when using the M062X functional), with the error asymptotic with the experimental measurements. Furthermore, the p K_a calculation results are independent of the reference

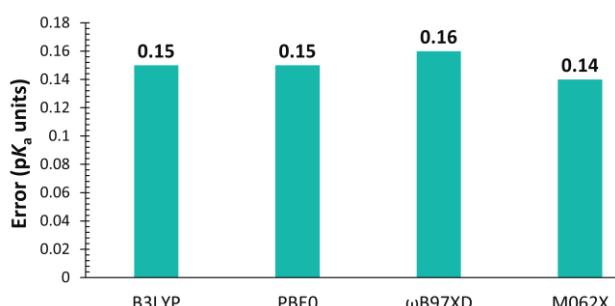


Figure 6. The mean absolute error (MAE) differences between pK_a calculated and experimental with statistical correction

compound. This is the recommended protocol for predicting pK_a values of phenolic compounds.

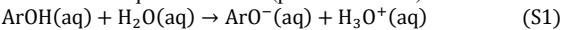
Acknowledgment: This research is funded by the University of Danang - Funds for Science and Technology Development under project number B2020-DN03-47.

REFERENCES

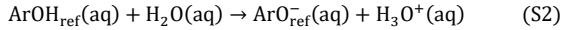
- [1] M. Namazian and S. Halvani, "Calculations of pK_a values of carboxylic acids in aqueous solution using density functional theory", *The Journal of Chemical Thermodynamics*, vol. 38, pp. 1495-1502, 2006.
- [2] K. Mansouri, N. F. Cariello, A. Korotcov, V. Tkachenko, C. M. Grulke, C. S. Sprankle, *et al.*, "Open-source QSAR models for pK_a prediction using multiple machine learning approaches", *Journal of cheminformatics*, vol. 11, pp. 1-20, 2019.
- [3] J. Jover, R. Bosque, and J. Sales, "QSPR prediction of pK_a for benzoic acids in different solvents", *QSAR & Combinatorial Science*, vol. 27, pp. 563-581, 2008.
- [4] G. C. Shields and P. G. Seybold, *Computational approaches for the prediction of pK_a values*: CRC Press, 2013.
- [5] P. Hunt, L. Hosseini-Gerami, T. Chrien, J. Plante, D. J. Ponting, and M. Segall, "Predicting pK_a Using a Combination of Semi-Empirical Quantum Mechanics and Radial Basis Function Methods", *Journal of chemical information and modeling*, vol. 60, pp. 2989-2997, 2020.
- [6] R. Casasnovas, J. Ortega-Castro, J. Frau, J. Donoso, and F. Munoz, "Theoretical pK_a calculations with continuum model solvents, alternative protocols to thermodynamic cycles", *International Journal of Quantum Chemistry*, vol. 114, pp. 1350-1363, 2014.
- [7] J. H. Jensen, C. J. Swain, and L. Olsen, "Prediction of pK_a Values for Druglike Molecules Using Semiempirical Quantum Chemical Methods", *The Journal of Physical Chemistry A*, vol. 121, pp. 699-707, 2017.
- [8] C. Andrés-Lacueva, A. Medina-Remón, R. Llorach, M. Urpi-Sarda, N. Khan, G. Chiva-Blanch, *et al.*, "Phenolic compounds: chemistry and occurrence in fruits and vegetables", ed: Wiley Online Library, 2010, pp. 53-80.
- [9] J. A. Vinson, Y. Hao, X. Su, and L. Zubik, "Phenol antioxidant quantity and quality in foods: vegetables", *Journal of agricultural and food chemistry*, vol. 46, pp. 3630-3634, 1998.
- [10] J. Chen, J. Yang, L. Ma, J. Li, N. Shahzad, and C. K. Kim, "Structure-antioxidant activity relationship of methoxy, phenolic hydroxyl, and carboxylic acid groups of phenolic acids", *Scientific reports*, vol. 10, pp. 1-9, 2020.
- [11] A. Aspée, C. Aliaga, L. Maretti, D. Zúñiga-Núñez, J. Godoy, E. Pino, *et al.*, "Reaction kinetics of phenolic antioxidants toward photoinduced pyranine free radicals in biological models", *The Journal of Physical Chemistry B*, vol. 121, pp. 6331-6340, 2017.
- [12] R. A. Albery, "Recommendations for nomenclature and tables in biochemical thermodynamics (IUPAC recommendations 1994)", *Pure and applied chemistry*, vol. 66, pp. 1641-1666, 1994.
- [13] S. Zhao, Z. Jin, and J. Wu, "New theoretical method for rapid prediction of solvation free energy in water", *The Journal of Physical Chemistry B*, vol. 115, pp. 6971-6975, 2011.
- [14] A. D. Becke, "A new mixing of Hartree-Fock and local density-functional theories", *The Journal of chemical physics*, vol. 98, pp. 1372-1377, 1993.
- [15] J. P. Perdew, M. Ernzerhof, and K. Burke, "Rationale for mixing exact exchange with density functional approximations", *The Journal of chemical physics*, vol. 105, pp. 9982-9985, 1996.
- [16] J.-D. Chai and M. Head-Gordon, "Long-range corrected hybrid density functionals with damped atom-atom dispersion corrections", *Physical Chemistry Chemical Physics*, vol. 10, pp. 6615-6620, 2008.
- [17] Y. Zhao and D. G. Truhlar, "The M06 suite of density functionals for main group thermochemistry, thermochemical kinetics, noncovalent interactions, excited states, and transition elements: two new functionals and systematic testing of four M06-class functionals and 12 other functionals", *Theoretical chemistry accounts*, vol. 120, pp. 215-241, 2008.
- [18] V. S. Bryantsev, "Predicting the stability of aprotic solvents in Li-air batteries: pK_a calculations of aliphatic C-H acids in dimethyl sulfoxide", *Chemical Physics Letters*, vol. 558, pp. 42-47, 2013.
- [19] F. Huang, J. Jiang, M. Wen, and Z.-X. Wang, "Assessing the performance of commonly used DFT functionals in studying the chemistry of frustrated Lewis pairs", *Journal of Theoretical and Computational Chemistry*, vol. 13, p. 1350074, 2014.
- [20] R. Krishnan, J. S. Binkley, R. Seeger, and J. A. Pople, "Self-consistent molecular orbital methods. XX. A basis set for correlated wave functions", *The Journal of chemical physics*, vol. 72, pp. 650-654, 1980.
- [21] S. Kheirjou, A. Abedin, A. Fattahi, and M. M. Hashemi, "Excellent response of the DFT study to the calculations of accurate relative pK_a value of different benzo-substituted quinuclidines", *Computational and Theoretical Chemistry*, vol. 1027, pp. 191-196, 2014.
- [22] F. M. Carvalho, Y. A. d. O. Só, A. S. K. Wernik, M. d. A. Silva, and R. Gargano, "Accurate acid dissociation constant (pK_a) calculation for the sulfachloropyridazine and similar molecules", *Journal of Molecular Modeling*, vol. 27, pp. 1-9, 2021.
- [23] S. Miertuš, E. Scrocco, and J. Tomasi, "Electrostatic interaction of a solute with a continuum. A direct utilization of AB initio molecular potentials for the revision of solvent effects", *Chemical Physics*, vol. 55, pp. 117-129, 1981.
- [24] A. Alibakhshi and B. Hartke, "Improved prediction of solvation free energies by machine-learning polarizable continuum solvation model", *Nature communications*, vol. 12, pp. 1-7, 2021.
- [25] M. Frisch, G. Trucks, H. Schlegel, G. Scuseria, M. Robb, J. Cheeseman, *et al.*, "Gaussian 16 revision a. 03. 2016; gaussian inc", Wallingford CT, vol. 2, 2016.
- [26] B. G. Tehan, E. J. Lloyd, M. G. Wong, W. R. Pitt, J. G. Montana, D. T. Manallack, *et al.*, "Estimation of pK_a using semiempirical molecular orbital methods. part 1: Application to phenols and carboxylic acids", *Quantitative Structure-Activity Relationships*, vol. 21, pp. 457-472, 2002.
- [27] M. D. Liptak, K. C. Gross, P. G. Seybold, S. Feldgus, and G. C. Shields, "Absolute pK_a determinations for substituted phenols", *Journal of the American Chemical Society*, vol. 124, pp. 6421-6427, 2002.
- [28] C. Book, "Chemical Book", [Online] Available: http://www.chemicalbook.com/ProductChemicalPropertiesCB8852597_EN.htm, 2017.

SUPPORTING INFORMATION (SI)

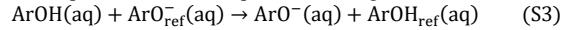
S1.1. The expression for calculating pK_a based on reference compound
The dissociation equation for acid (phenolic: ArOH) in solution:



The dissociation equation for the reference phenolic (ArOH_{ref}) in solution:



Subtract Equation (S2) from Equation (S1) to get:



Since Gibbs energy is a state function:

$$\Delta G_T^0(\text{S3}) = \Delta G_T^0(\text{S1}) - \Delta G_T^0(\text{S2}) \quad (\text{S4})$$

Where, $\Delta G_T^0(\text{Si})$ is the standard Gibbs energy of the reactions in terms of Eq (Si). The expression of standard Gibbs energy in terms of acid dissociation constant:

$$\Delta G_T^0(\text{S1}) = -RT \ln K_a \quad (\text{S5})$$

$$\Delta G_T^0(\text{S3}) = -RT \ln K_a^{ref} \quad (\text{S6})$$

Where, K_a and K_a^{ref} are the acid dissociation constant of ArOH and ArOH_{ref} . Combining equations (S4), (S5) and (S6) get:

$$\Delta G_T^0(\text{S3}) = -RT \ln \frac{K_a}{K_a^{ref}} \quad (\text{S7})$$

Equation (S7) can be rewritten:

$$-\log K_a = \frac{\Delta G_T^0(3)}{RT \ln 10} - \log K_a^{ref} \quad (\text{S8})$$

To the definition $pK_a = -\log K_a$, finally get the expression pK_a in terms of pK_a^{ref} :

$$pK_a = \frac{\Delta G_T^0(5)}{RT \ln 10} + pK_a^{ref} \quad (\text{S9})$$

