

Title	Developing a design guideline of boronic acid derivatives to scavenge targeted sugars in the formose reaction products using DFT-based machine learning
Author(s)	Ishihara, Nanako; Chikatani, Genta; Nishijima, Hiroaki et al.
Citation	Chemistry Letters. 2024, 53(6), p. upae087
Version Type	АМ
URL	https://hdl.handle.net/11094/97753
rights	
Note	

The University of Osaka Institutional Knowledge Archive : OUKA

https://ir.library.osaka-u.ac.jp/

The University of Osaka

# Developing a Design Guideline of Boronic Acid Derivatives to Scavenge Targeted Sugars in the Formose Reaction Products using DFT-based Machine Learning

Nanako Ishihara<sup>1</sup>, Genta Chikatani<sup>1</sup>, Hiroaki Nishijima<sup>1</sup>, Hiro Tabata<sup>1</sup>, Yoko Hase<sup>1</sup>, Yoshiharu Mukouyama<sup>1,2</sup>, Shuji Nakanishi<sup>1,3,\*</sup>, Shiho Mukaida<sup>1,\*</sup>

<sup>1</sup> Research Center for Solar Energy Chemistry, Graduate School of Engineering Science, Osaka University, 1-3 Machikaneyama, Toyonaka, Osaka 560-8531, Japan

<sup>2</sup> Division of Science, College of Science and Engineering, Tokyo Denki University, Hatoyama, Saitama 350-0394, Japan

<sup>3</sup> Innovative Catalysis Science Division, Institute for Open and Transdisciplinary Research Initiatives (ICS-OTRI), Osaka University, Suita, Osaka 565-0871, Japan

\*Corresponding author: Research Center for Solar Energy Chemistry, Graduate School of Engineering Science, Osaka University, 1-3 Machikaneyama, Toyonaka, Osaka 560-8531, Japan Email: <u>nakanishi.shuji.es@osaka-u.ac.jp</u> (S.Nakanishi), <u>mukaida@rcsec.chem.es.osaka-u.ac.jp</u> (S.Mukaida)

# Abstract

Formose reaction facilitates the synthesis of sugars from HCHO, yet the valuable sugars constitute only a small portion of the total products. This necessitates the need for a chemical scavenger capable of selectively capturing only valuable sugars. With over 600,000 potential combinations of boronic acid-based scavengers available, pursuing a deductive search approach is unfeasible. This study aims to derive guidelines for designing scavengers that readily bind with target sugars while avoiding non-target ones, via machine learning informed by DFT calculations.

Keywords: Formose reaction, Machine learning, DFT calculation

## Graphical abstract



Since its discovery by Butlerov et al. in 1861,<sup>1</sup> the formose 2 reaction-a non-enzymatic pathway for synthesizing sugars from formaldehyde (HCHO) under basic 3 4 conditions-has attracted considerable interest for its 5 implications in origin of life and food production.<sup>2-6</sup> Recently, our research successfully demonstrated that 6 7 Corvnebacterium glutamicum can be cultured using formose 8 sugar, highlighting its potential as a substrate for the bioproduction of valuable compounds.<sup>7</sup> The formose reaction 9 involves a myriad of reaction combinations featuring 10 11 carbonyl groups, including aldol reactions, retro-aldol 12 reactions, and aldose-ketose isomerization reactions, thereby

In this study, we clarified the characteristics that a selective scavenger for bio-assimilable sugars contained in formose reaction products must possess, utilizing machine learning informed by DFT calculations.

13 forming a highly complex chemical reaction network.<sup>8</sup>
14 Therefore, achieving high selectivity for desired products
15 within this network is challenging.

16 Developing selective catalysts has been proposed as a 17 strategy to address the challenge inherent in the formose reaction.9-14 As an approach different from catalyst 18 19 development, a strategy has also been proposed that involves 20 adding scavengers that bind to the target sugar. Through the 21 protection of hydroxy groups in these target sugars by the 22 chemical scavenger, further reactions of the target sugars can 23 be prevented, thereby increasing their yields. Boronic acids 24 (BA) emerge as promising materials for these scavengers, as

they can form boronic esters with the diols of sugars in 1 2 aqueous solutions, thereby stabilizing the products.<sup>15</sup> In fact, Ricardo et al. have discovered that boronic acid, protect 3 ribose in the formose reaction system.<sup>16</sup> However, due to their 4 5 general property of forming cyclic esters with diols, BA can 6 bind not only to the target sugars in the formose reaction 7 products but also to non-target sugars, which are unnecessary. 8 Thus, simply adding BA to the formose reaction system 9 captures non-target sugars as well, potentially clogging the 10 entire reaction.17

11 This issue can be addressed by employing boronic acid derivatives that have strong binding affinity for the target 12 sugars while having weak binding affinity for the non-target 13 14 sugars. To find suitable boronic acid derivatives for the 15 formose reaction, exploration of all combination of boronic 16 acid derivatives and sugars is needed. However, there are 17 more than 3,000 types of commercially available boronic acid 18 derivatives. Additionally, the formose reaction can yield up 19 to 23 types of sugars, not even accounting for stereoisomers, 20 presenting a significant challenge in identifying the most 21 effective derivatives for selective sugar capture. Moreover, 22 given that a single sugar molecule can have multiple pairs of 23 hydroxyl groups capable of interacting with boronic acid 24 derivatives<sup>18</sup>, there is more than one potential interaction 25 pattern between a sugar and a boronic acid derivatives; in 26 some cases, there can be as many as eight varieties. 27 Consequently, even when restricting the variety of sugars 28 under consideration, the total number of potential 29 combinations can surpass 600,000. Exploring these 30 combinations through deductive methods is unfeasible. In addressing these challenges, employing machine learning proves to be a valuable approach.<sup>19-23</sup> Hence, this study aims 31 32 33 to derive guidelines for designing scavengers that enhance 34 selectivity for target sugars by investigating boronic acid 35 derivative-based scavengers through machine learning, which is informed by density functional theory (DFT) 36

calculations. There have been several examples where
machine learning has been applied to formose reaction
systems<sup>24,25</sup>; however, this study is the first to explore
scavenger materials using this approach.

41 Figure 1(a) shows an overview of the formose reaction 42 network, highlighting both the target and non-target sugars 43 investigated in this study. All sugar structure used in this 44 work is detailed in Table S1. In this study, bio-assimilable 45 sugars such as glucose, fructose, ribose, and arabinose, which 46 possess five (C5) and six (C6) carbon atoms, were defined as 47 target sugars. On the other hand, aldoses (C4a) and ketoses 48 (C4k), containing four carbon atoms, were classified as non-49 target sugars. This classification stems from the fact that C4a 50 and C4k serve as precursors to C5 and C6 sugars, and their 51 interception by scavengers might inhibit the generation of the 52 larger C5 and C6 sugars. Furthermore, the reaction products 53 of C4k and C4a with glycolaldehyde (C2) (C4k+C2, 54 C4a+C2), the compound formed by the reaction of C4k with 55 two molecules of HCHO (C1) (C4k+C1+C1), and the 56 product resulting from the combination of two molecules of 57 C3k (C3k+C3k) were designated as non-target sugars in this 58 study. This is due to their potential to transform into glucose 59 or fructose through subsequent reactions within the formose 60 reaction network.

61 Figure 1(b) depicts the conceptual framework of this 62 study. We computed the Gibbs free energy change ( $\Delta_r G^\circ$ ) for 2,927 reactions from the possible 615,876 combinations of 63 ester formation reactions between boronic acid derivatives 64 and sugars using DFT calculations. Utilizing these 65 computations as training data, we developed a machine 66 learning model to predict the  $\Delta_r G^{\circ}$  for the entire set of 67 68 615,87616 ester formation reactions. Given the aim of this 69 study to identify boronic acid derivatives that selectively bind 70 to target sugars, our focus is on the difference between the 71 average  $\Delta_r G^{\circ}$  for ester formation reactions involving a 72 specific boronic acid with the 4 types of target sugars



Figure 1. (a) The formose reaction network. Species surrounded with red and blue squares represents target sugars and non-target sugars, respectively. (b) The conceptual framework of this study and the actual versus predicted plot from the constructed machine learning model.

 $(\overline{\Delta_r G^{\circ}}_{target})$  and the average  $\Delta_r G^{\circ}$  for ester formation 1 2 reactions between that same boronic acid and the 7 types of 3 non-target sugars ( $\overline{\Delta_r G^{\circ}}_{non-target}$ ). This difference is represented as  $\Delta_r G_{diff}$ . 4 5

 $\Delta_r G_{\text{diff}} = \overline{\Delta_r G^{\circ}}_{\text{non-target}} - \overline{\Delta_r G^{\circ}}_{\text{target}} (1)$ 

6 The greater the  $\Delta_r G_{\text{diff}}$  value, the more it exhibits 7 characteristics that align with the objective of this study.  $\Delta_r G_{\text{diff}}$  values are derived from the  $\Delta_r G^\circ$  values predicted by 8 9 the machine learning model, and boronic acid derivatives that 10 appear promising are identified based on these values.

11 Finally, the  $\Delta_r G^{\circ}$  values for the ester formation reactions 12 between the boronic acid derivative identified as optimal by 13 the machine learning model and both target and non-target 14 sugars are calculated using DFT. These values are then 15 compared with those of a standard boronic acid (i.e., phenylboronic acid (PBA)) to validate the suitability of the 16 17 boronic acid derivative proposed by the machine learning 18 model. We would like to re-emphasize here that the purpose 19 of this study is to perform a first screening to extract the 20 general trends that promising boronic acid derivatives should 21 possess, through comprehensive analysis of over 600,000 22 combinations. Therefore, the accuracy of the machine 23 learning and DFT calculations presented hereafter, as well as 24 the interpretation of their results, are conducted in light of this 25 objective.

26 The actual versus predicted plot for the constructed 27 machine learning model is displayed in Figure 1b. The 28 coefficient of determination (R<sup>2</sup>) exceeded 0.8, as indicated 29 in Figure 1(b), verifying that the model has adequate performance. The results of the analysis of feature 30 importance by SHAP values<sup>26</sup> are shown in Figure 2(a). 31 32 Notably, among the different factors considered, the sugar 33 species emerged as the most significant feature, suggesting 34 that the structure of the sugar plays a predominant role in 35 determining the  $\Delta_r G^\circ$ , while the contribution from boronic 36 acid functional groups is relatively minor. This finding is consistent with the discussions in the field of glucose sensing, 37 38 where it has been argued that the sugar species, more 39 specifically its structure, significantly influences its affinity with boronic acid derivatives.<sup>16</sup> The fact that the sugar species, 40 41 rather than the boronic acid derivatives or the boronic ester 42 formation products, was identified as the most significant 43 factor in the  $\Delta_r G^\circ$  calculation further supports the validity of the machine learning approach employed in this study.<sup>27,28</sup> 44

45 The findings from the SHAP analysis further revealed 46 that the descriptor Chi2n, indicative of the topological 47 characteristics of molecules, also played a significant role. 48 The Chi2n value, understood to be calculated by considering 49 interatomic paths, increases in response to greater structural 50 complexity, such as branching, and with the enlargement of 51 molecular size (Table S2).<sup>29</sup> The positive correlation between 52 the Chi2n values and the predicted  $\Delta_r G^{\circ}$  values, as confirmed 53 by the SHAP analysis depicted in Figure 2(b), suggests that 54 steric hindrance is a primary contributing factor. Additionally, 55 the analysis revealed that descriptors associated with the C=N 56 bond in boronic acid derivatives, specifically PEOE VSA9 57 (Figure S1), and VSA ESTATE6, which represents the 58 carbon of the phenyl group (Figure S2), did not significantly 59 contribute. Taken all together, it is suggested that the steric



Figure 2. (a) The SHAP importance values of the machine learning model constructed. Blue, orange, and green represents sugars, boronic acid derivatives, and ester bond formation products, respectively. (b) A dependence plot between Chi2n and SHAP values.

60 effects associated with boronic acid derivatives play a more 61 substantial role in influencing the reaction  $\Delta_r G^{\circ}$  than the 62 presence of specific functional groups or bonds.

63 As previously discussed, boronic acid derivatives were 64 assessed based on their  $\Delta_r G_{\text{diff}}$  value (equation (1)). In 65 comparison to the  $\Delta_r G_{diff}$  value of 3.344, associated with standard PBA, the  $\Delta_r G_{diff}$  value for the boronic acid 66 derivative Bortezomib-identified as an optimal molecule by 67 68 machine learning-was significantly higher, registering at 5.913, as detailed in Table S3. Thus, through DFT 69 70 calculations, it has been confirmed that Bortezomib, as 71 suggested by the machine learning model, indeed exhibits 72 better selectivity towards target sugars compared to PBA.

73 The increase in the  $\Delta_r G_{\text{diff}}$  value for Bortezomib may 74 result from two scenarios: (1) higher affinity for the target 75 sugar, or (2) lower affinity for the non-target sugars. 76 Consequently, we investigated which scenario, (1) or (2), 77 applies. The  $\Delta_r G^{\circ}$  values, as estimated by DFT calculations, for the ester formation reactions involving Bortezomib and 78 79 PBA with their target and non-target sugars are shown in 80 Figure 3a and 3b, respectively. It should be noted here that we plot the average  $\Delta_r G^{\circ}$  values of the ester formation 81



1 reactions between various isomers of each sugar and the

Figure 3. The average reaction Gibbs energy of (a) Bortezomib and (b) PBA for each sugar, and (c) the difference in reaction Gibbs energy between Bortezomib and PBA.

2 boronic acid derivatives. For the target sugars, the values 3 were around -30 kcal/mol with Bortezomib and -40 kcal/mol 4 with PBA. In contrast, for the non-target sugars, the values 5 exceeded -20 kcal/mol with Bortezomib and were around -30 6 kcal/mol with PBA. The bar plot of Figure 3c illustrates the difference in  $\Delta_r G^\circ$  values between the two boronic acid 7 8 derivatives. On the other hand, the dotted lines in Figure 3c 9 represent the average differences in the  $\Delta_r G^{\circ}$  values between 10 Bortezomib and PBA for each target and non-target sugars. 11 (That is, they do not correspond to the simple average of each 12 bar shown in Fig. 3c.) This average yields positive values for 13 both target and non-target sugars, indicating that the larger 14  $\Delta_r G_{\text{diff}}$  value in Bortezomib, compared to PBA, was 15 attributable to reason (2) rather than (1).

Next, we will explore the chemical reasons behind 16 the increased difference in the  $\Delta_r G_{diff}$  value. As shown in 17 18 Figure 3(c), focusing on the magnitude of the  $\Delta_r G^{\circ}$  for each 19 sugar with PBA and Bortezomib, it is apparent that 20 Bortezomib tends to exhibit larger  $\Delta_r G^\circ$  compared to PBA. 21 Given the discussion regarding Chi2n in Figure 2, the 22 primary reason for this is believed to be the steric hindrance 23 resulting from the molecular complexity of Bortezomib. On the other hand, the aforementioned increase in  $\Delta_r G^{\circ}$  in the 24 25 Bortezomib varies depending on the sugar species, being 26 smaller for target sugars compared to non-target sugars. This 27 variation is thought to stem from the inherent structure of the 28 sugars themselves. Specifically, cyclic fructose and cyclic 29 ribose in target sugars have many cis-configured OH groups 30 advantageous for the formation of boronic acid esters, and the 31 sugar itself has small steric hindrance.<sup>30</sup> Conversely, the 32 branched sugars contained in the non-target sugars. In other 33 words, in the target sugars, the  $\Delta_r G^{\circ}$  mentioned earlier is 34 mitigated, and such an effect cannot be expected in the non-35 target sugar. Therefore, it is speculated that in Bortezomib, 36 the difference in  $\Delta_r G^\circ$  between the target sugars and non-37 target sugars becomes significant, suggesting its potential 38 function as a selective scavenger for sugars.

39 In this study, we aimed to obtain a general guideline for 40 designing boronic acid-based scavengers that are efficacious 41 in selectively capturing target sugars. From a comprehensive 42 pool of 615,876 possible combinations, involving boronic 43 ester products derived from 42 types of monosaccharides and 44 3003 commercially available monoboronic acid molecules, 45 Bortezomib emerged as an optimal candidate. The insights from SHAP analysis and DFT calculations indicate that 46 47 controlling the spatial structure is more crucial than tuning 48 electronic properties, such as the selection of functional 49 groups, in the context of boronic acid derivatives. In this 50 study, we searched for the ideal scavenger under the 51 assumption that one molecule of a boronic acid derivative 52 forms an ester bond with one sugar molecule, as a first 53 screening. Through machine learning, Bortezomib was 54 proposed as a candidate molecule. This study has established 55 a fundamental guideline demonstrating that effective control 56 of stereo-configuration is beneficial. Following this guideline, 57 future work will likely find that enhancing reaction 58 selectivity can be effectively achieved by combining multiple 59 mono-boronic acid molecules or poly-boronic acids to a 60 target sugar.

#### Supplementary data 62

61

64

63 Supplementary material is available at Chemistry Letters

#### 65 Acknowledgements

66 This work was achieved through the use of SQUID at the 67 Cybermedia Center, Osaka University. I would like to thank 68 Kaito Nagita for advice about conducting DFT calculations 69 on SQUID and useful discussions. 70

#### 71 References

- 1 A. Butlerov, Justus Liebigs Ann. Chem. 1861, 120, 295.
- 2 N. W. Gabel, C. Ponnamperuma, Nature 1967, 216, 453-455.
- 3 C. Reid, L. E. Orgel, Nature 1967, 216, 455.

J. B. García Martínez, K. A. Alvarado, X. Christodoulou, D. 4 C. Denkenberger, J. CO2 Util. 2021, 53, 101726.

72 73 74 75 76 77 78 79 S. Cestellos-Blanco, S. Louisia, M. B. Ross, Y. Li, N. E. 5 Soland, T. C. Detomasi, J. N. Cestellos Spradlin, D. K. Nomura, P. Yang, Joule 2022, 6, 2304.

80 6 I. V. Delidovich, A. N. Simonov, O. P. Taran, V. N. Parmon, 81 ChemSusChem 2014, 7, 1833.

82 7 H. Tabata, H. Nishijima, Y. Yamada, R. Miyake, K.

83 Yamamoto, S. Kato, S. Nakanishi, ChemBioChem 2024, 25, e202300760.

84 W. E. Robinson, E. Daines, P. van Duppen, T. de Jong, W. T. 8

85 S. Huck, Nature Chem. 2022, 14, 623.

1

- 9 J. Castells, F. López-Calahorra, F. Geijo, Carbohydr. Res. 1983, 116, 197.
- 23 H. Tabata, G. Chikatani, H. Nishijima, T. Harada, R. Miyake, 10 S. Kato, K. Igarashi, Y. Mukouyama, S. Shirai, M. Waki, Y. Hase, S.
- 4 5 Nakanishi, Chem. Sci. 2023, 14, 13475.
- 6 7 8 9 M. Waki, S. Shirai, Y. Hase, Dalton Transactions 2024, 53, 11 2678.
- 12 Z. Iqbal, S. Novalin, Curr. Org. Chem. 2012, 16, 769.
- 13 T. I. Khomenko, M. M. Sakharov, O. A. Golovina, Russ. 10 Chem. Rev. 1980, 49, 570.
- 11 14 T. Matsumoto, H. Yamamoto, S. Inoue, J. Am. Chem. Soc. 12 1984, 106, 4829.
- 13 15 J. P. Lor, J. O. Edwards, J. Org. Chem. 1959, 24, 769.
- 14 15 16 A. Ricardo, M. A. Carrigan, A. N. Olcott, S. A. Benner, Science 2004, 303, 196.
- 16 T. Imai, T. Michitaka, A. Hashidzume, Beilstein J. Org. 17 17 Chem. 2016, 12, 2668.
- 18 T. D. James, K. R. A. Samankumara Sandanayake, S. Shinkai, 18 Angew Chem. Int. Ed. Engl. 1996, 35, 1910.
- Y. Guo, X. He, Y. Su, Y. Dai, M. Xie, S. Yang, J. Chen, K. 19 Wang, D. Zhou, C. Wang, J. Am. Chem. Soc. 2021, 143, 5755.
- D. P. Metcalf, Z. L. Glick, A. Bortolato, A. Jiang, D. L. 20 Cheney, C. D. Sherrill, J. Chem. Inf. Model. 2023, 64.
- 21 T. Ando, N. Shimizu, N. Yamamoto, N. N. Matsuzawa, H.
- Maeshima, H. Kaneko, J. Phys. Chem. A 2022, 126, 6336.
- S. Samizo, H. Kaneko, ACS Omega 2023, 8, 27247. 22
- 23 A. Yada, K. Nagata, Y. Ando, T. Matsumura, S. Ichinoseki, K. Sato, Chem. Lett. 2018, 47, 284.
- H. J. Cleaves, G. Hystad, A. Prabhu, M. L. Wong, G. D. Cody, 24
- S. Economon, R. M. Hazen, Proc. Natl. Acad. Sci. USA 2023, 120, e2307149120.
- 25 S. Asche, G. J. T. Cooper, G. Keenan, C. Mathis, L. Cronin, Nat. Commun. 2021, 12, 1.
- 26 S. M. Lundberg, P. G. Allen, S.-I. Lee, Adv. Neural Inf. Process. Syst. 2017, 10, 4768-4777.
- 27 H. Fang, G. Kaur, B. Wang, J. Fluoresc. 2004, 14, 481.
- 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 45 36 37 38 940 28 N. Fujita, S. Shinkai, T. D. James, Chem. Asian J. 2008, 3, 1076.
  - 29 L. H. Hall, L. B. Kier, Rev. Comput. Chem. 1991, 367.
  - 30 Y. Suzuki, M. Shimizu, T. Okamoto, T. Sugaya, S. Iwatsuki,
- 41 42 M. Inamo, H. D. Takagi, A. Odani, K. Ishihara, ChemistrySelect 2016, 1, 5141.
- 43