Supplementary Information

# Molecular screening for solid–solid phase transition by machine learning

Daisuke Takagi,[1] Kazuki Ishizaki,[2] Toru Asahi,[1,2] Takuya Taniguchi*[3]

[1] Department of Life Science and Medical Bioscience, Graduate School of Advanced Science and Engineering, Waseda University, 3-4-1 Okubo, Shinjuku-ku, Tokyo 169-8555, Japan

[2] Department of Advanced Science and Engineering, Graduate School of Advanced Science and Engineering, Waseda University, 3-4-1 Okubo, Shinjuku-ku, Tokyo 169-8555, Japan

[3] Center for Data Science, Waseda University, 1-6-1 Nishiwaseda, Shinjuku-ku, Tokyo 169-8050, Japan

* Correspondence to takuya.taniguchi@aoni.waseda.jp

**<u>Contents</u>**

# Supplementary Figures



**Figure S1. Molecules collected as positive data.**

(a)



| | $T_{endo}$ | $T_{endo(max)}$ | $T_{endo(min)}$ |
|---|---|---|---|
| Count | 144 | 88 | 88 |
| Mean | 343.5 | 350.6 | 324.2 |
| Std | 95.0 | 95.5 | 101.0 |
| Min | 80.0 | 80.0 | 80.0 |
| Max | 533.0 | 533.0 | 518.0 |

(b)



| | $T_{exo}$ | $T_{exo(max)}$ | $T_{exo(min)}$ |
|---|---|---|---|
| Count | 110 | 72 | 72 |
| Mean | 308.0 | 314.3 | 291.4 |
| Std | 93.0 | 92.1 | 96.3 |
| Min | 120.0 | 135.0 | 120.0 |
| Max | 564.2 | 564.2 | 463.0 |

(c)



| | $\Delta H_{endo}$ | $\Delta H_{endo(max)}$ | $\Delta H_{endo(min)}$ |
|---|---|---|---|
| Count | 64 | 46 | 46 |
| Mean | 4.3 | 4.8 | 3.9 |
| Std | 6.2 | 6.8 | 6.1 |
| Min | 0.04 | 0.04 | 0.04 |
| Max | 36.4 | 36.4 | 36.4 |

(d)



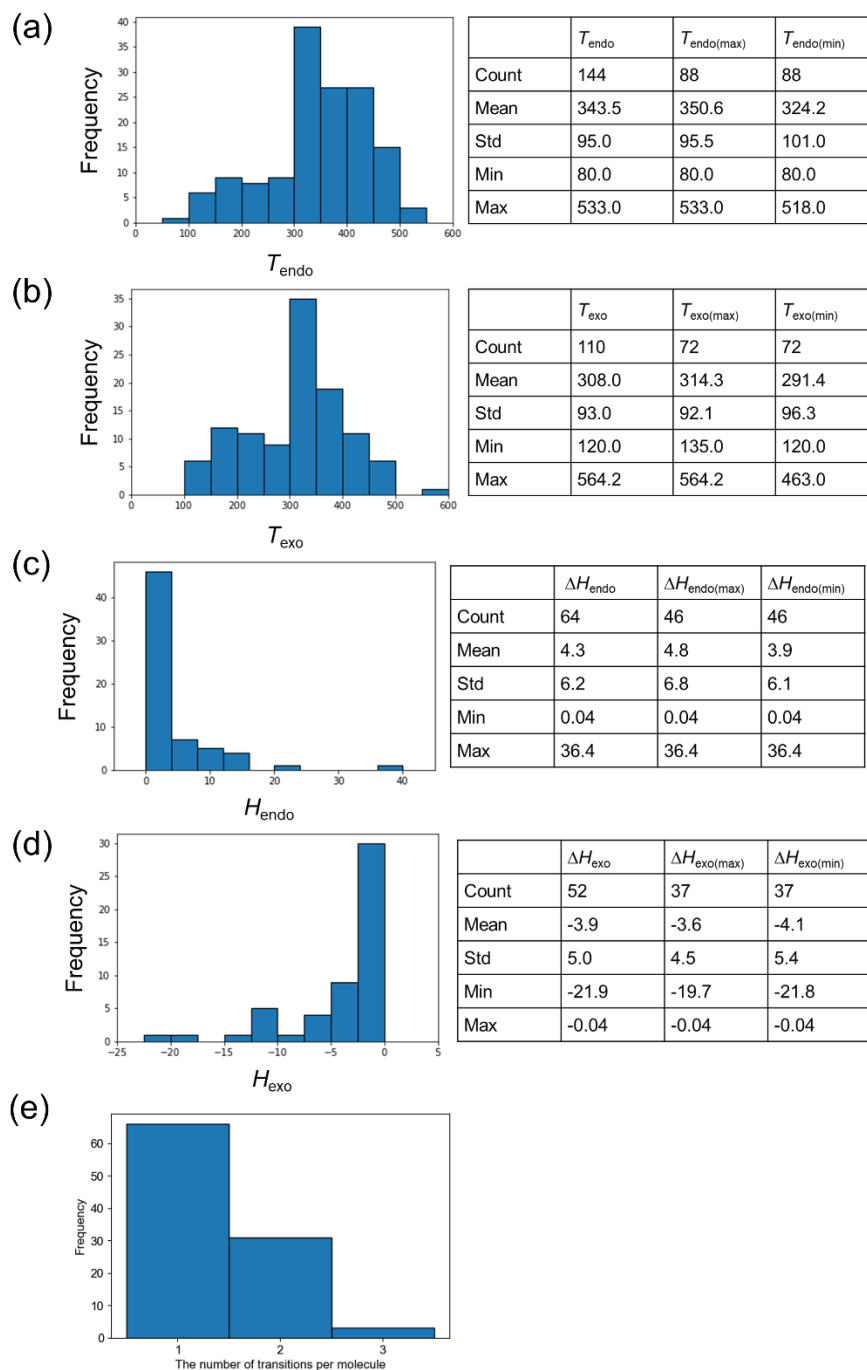| | $\Delta H_{exo}$ | $\Delta H_{exo(max)}$ | $\Delta H_{exo(min)}$ |
|---|---|---|---|
| Count | 52 | 37 | 37 |
| Mean | -3.9 | -3.6 | -4.1 |
| Std | 5.0 | 4.5 | 5.4 |
| Min | -21.9 | -19.7 | -21.8 |
| Max | -0.04 | -0.04 | -0.04 |

(e)



**Figure S2. The distribution of positive dataset.** (a,b) Statistics of transition temperature corresponding to (a) endothermic and (b) exothermic transitions. (c,d) Statistics of transition enthalpy corresponding to (c) endothermic and (d) exothermic transitions. (e) The distribution of the number of solid phase transitions per molecule.
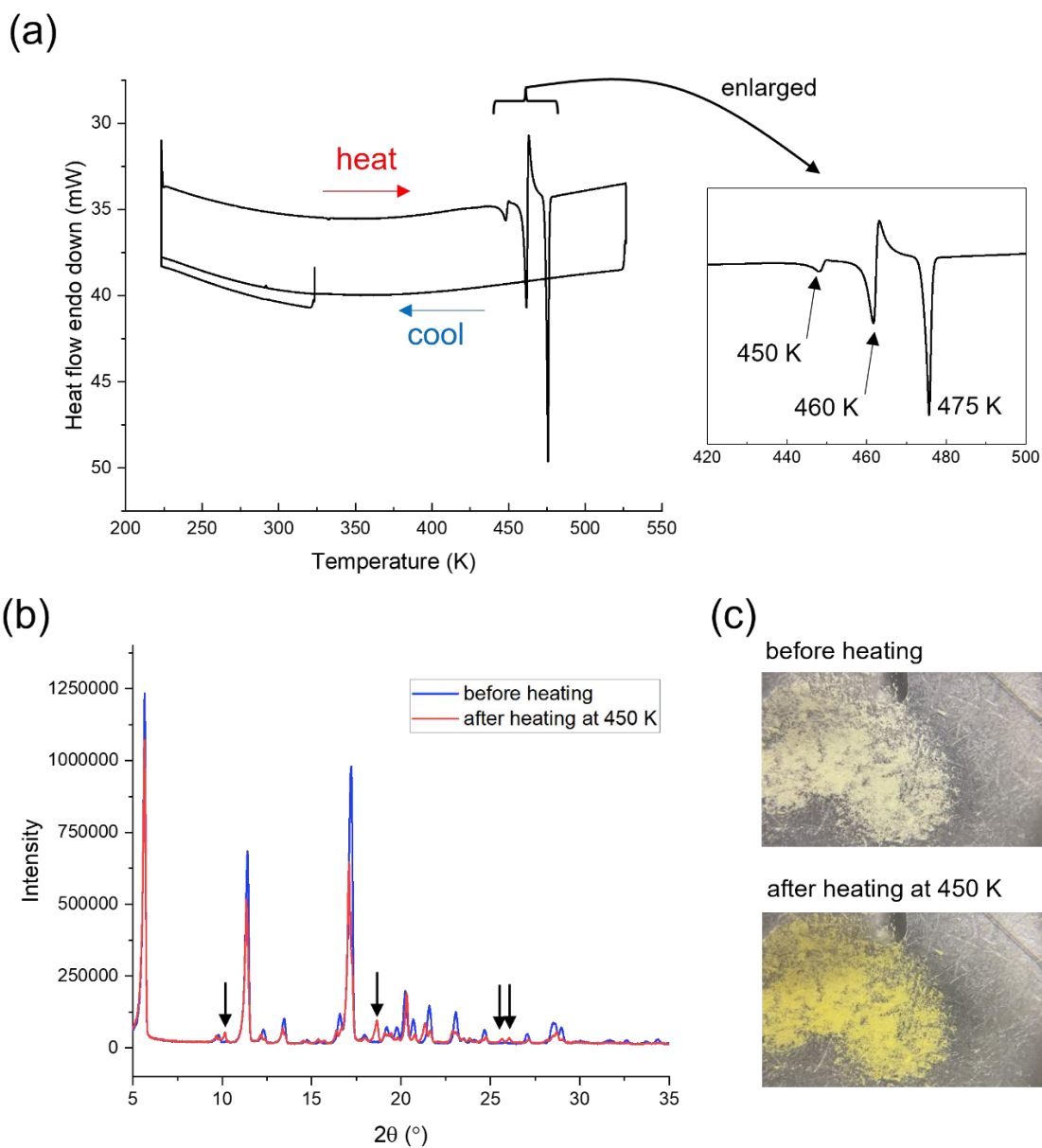
**Figure S3. Crystal-to-crystal phase transition of the crystal of OCAPAK.** (a) Thermal analysis by DSC. The crystal exhibited a crystal-to-crystal transition at 450 K, and an anomalous endothermic-exothermic peak at 460 K before melting at 475 K. (b) XRD patterns of the crystal before heating, and after heating at 450 K. Both measurements were performed at 293 K. Unique peaks after heating are indicated by arrows. The slight difference of XRD pattern before and after heating reflects the crystal-to-crystal phase transition. (c) Pictures of crystalline powder before and after heating. Crystal appearance did not change before and after heating, although slight color change was observed.

**Figure S4. Molecular structures suggested by PU learning with *p* = 0.2.** References are summarized in supplementary references.

**Supplementary Tables**

**Table S1. Hyperparameters used for ML models in classification and regression.**

| ML models | Hyperparameters |
|---|---|
| Classification | |
| RF | n_estimaters:100, max_depth: None |
| NN | hidden_layer_sizes: (50,), activation: 'relu' |
| SVM | kernel: 'rbf', gamma: 0.2 |
| GBDT | n_estimaters:100, max_depth: 3 |
| Regression | |
| NN | hidden_layer_sizes: (100,), activation: 'relu' |
| RF | n_estimaters: 100, max_depth: None |
| TL-NN | (Search space) |
| | n_layers: {1, 10}, n_dim: {50, 100, 200, 300, 400, 500, 750, 1000} |
| ECFP | n_layers: 5, n_dim: 1000 |
| Avalon | n_layers: 4, n_dim: 750 |
| ErG | n_layers: 3, n_dim: 1000 |
| MACCSKeys | n_layers: 4, n_dim: 1000 |
| Estate | n_layers: 7, n_dim: 1000 |

**Table S2. Comparison of true positive rate (TPR) between PU and BC.** The value is the averaged TPR in 10-fold CV.

| | PU | | | | BC | | | |
|---|---|---|---|---|---|---|---|---|
| | RF | NN | SVM | GBDT | RF | NN | SVM | GBDT |
| Mordred | 0.13 | 0.0 | 0.30 | 0.96 | 0.06 | 0.0 | 0.0 | 0.02 |
| ECFP | 0.22 | 0.27 | 0.39 | 0.62 | 0.14 | 0.13 | 0.01 | 0.06 |
| Avalon | 0.28 | 0.21 | 0.52 | 0.28 | 0.14 | 0.19 | 0.0 | 0.01 |
| ErG | 0.26 | 0.21 | 0.05 | 0.46 | 0.12 | 0.11 | 0.0 | 0.0 |
| RDKit | 0.20 | 0.01 | 0.0 | 0.74 | 0.06 | 0.0 | 0.0 | 0.01 |
| MACCSKeys | 0.15 | 0.21 | 0.20 | 0.04 | 0.10 | 0.14 | 0.0 | 0.0 |
| Estate | 0.16 | 0.03 | 0.0 | 0.14 | 0.10 | 0.0 | 0.0 | 0.0 |

**Table S3. Comparison of TPR and SE between each combination of molecular descriptor and classification model.** Each value is the average obtained from 10-fold CV.

| | TPR | | | | SE | | | |
|---|---|---|---|---|---|---|---|---|
| | RF | NN | SVM | GBDT | RF | NN | SVM | GBDT |
| Mordred | 0.13 | 0.0 | 0.30 | 0.96 | 69.2 | 4936.1 | 1.0 | 1.0 |
| ECFP | 0.22 | 0.27 | 0.39 | 0.62 | 87.8 | 1110.4 | 12441.4 | 2.1 |
| Avalon | 0.28 | 0.21 | 0.52 | 0.28 | 90.6 | 2011.0 | 21990.2 | 113.3 |
| ErG | 0.26 | 0.21 | 0.05 | 0.46 | 72.9 | 384.8 | 64422.4 | 73.3 |
| RDKit | 0.20 | 0.01 | 0.0 | 0.74 | 248.8 | 8575.1 | NaN* | 1.0 |
| MACCSKeys | 0.15 | 0.21 | 0.20 | 0.04 | 67.0 | 333.7 | 13594.6 | 145.3 |
| Estate | 0.16 | 0.03 | 0.0 | 0.14 | 74.0 | 202.2 | 171825.4 | 111.6 |

* SE was not able to be calculated because $n_{up}$ was zero in the combination of RDKit-SVM.

**Table S4. Regression performance based on MAE when using molecular descriptors for each property.** Each MAE value represents the mean and standard deviation on five 5-fold CVs. The letter "-" in TL-NN means the transfer wasn't succeeded due to the mismatch of input dimension. TL-NN is the best model among all the fine-tuning models.

| | NN | TL-NN | RF | mean |
|---|---|---|---|---|
| $T_{endo(max)}$ (K) | | | | |
| Mordred | $2.0 \times 10^7$ ($2.4 \times 10^7$) | - | 58.7 (8.1) | 75.4 |
| ECFP | 227.8 (21.1) | 76.7 (15.8) | 71.6 (9.9) | 75.4 |
| Avalon | 157.1 (25.0) | 76.2 (11.8) | 64.1 (12.5) | 75.4 |
| ErG | 206.8 (19.4) | 94.1 (21.7) | 66.0 (13.2) | 75.4 |
| RDkitDesc | 149.1 (34.4) | - | 70.6 (13.1) | 75.4 |
| MACCSKeys | 183.9 (29.9) | 77.9 (10.6) | 69.5 (14.3) | 75.4 |
| Estate | 220.4 (25.3) | 98.2 (20.5) | 69.9 (11.7) | 75.4 |
| $T_{endo(min)}$ (K) | | | | |
| Mordred | $2.4 \times 10^7$ ($3.1 \times 10^7$) | - | 64.9 (9.0) | 80.3 |
| ECFP | 200.2 (24.7) | 75.3 (9.8) | 69.8 (14.3) | 80.3 |
| Avalon | 143.1 (19.8) | 75.2 (17.0) | 71.3 (14.4) | 80.3 |
| ErG | 188.6 (28.8) | 98.9 (24.1) | 68.0 (14.2) | 80.3 |
| RDkitDesc | 147.5 (30.4) | - | 71.1 (17.4) | 80.3 |
| MACCSKeys | 167.9 (26.5) | 80.9 (13.0) | 69.9 (12.1) | 80.3 |
| Estate | 197.9 (20.3) | 96.3 (18.0) | 74.3 (17.5) | 80.3 |
| $T_{exo(max)}$ (K) | | | | |
| Mordred | $1.8 \times 10^7$ ($2.5 \times 10^7$) | - | 66.2 (14.4) | 71.5 |
| ECFP | 196.9 (24.9) | 78.3 (16.6) | 70.5 (14.0) | 71.5 |
| Avalon | 151.7 (21.0) | 76.1 (19.1) | 65.0 (16.4) | 71.5 |
| ErG | 184.1 (31.0) | 105.2 (28.0) | 65.7 (16.2) | 71.5 |
| RDkitDesc | 132.6 (43.7) | - | 69.8 (17.3) | 71.5 |
| MACCSKeys | 161.5 (26.1) | 81.8 (14.2) | 72.9 (11.6) | 71.5 |
| Estate | 187.4 (27.6) | 100.5 (20.1) | 69.7 (12.9) | 71.5 |

| $T_{\text{exo(min)}}$ (K) | | | | |
|---|---|---|---|---|
| Mordred | $3.5\times10^7$ ($3.7\times10^7$) | - | 69.4 (12.8) | 80.8 |
| ECFP | 173.0 (22.7) | 76.2 (16.2) | 70.8 (13.7) | 80.8 |
| Avalon | 137.4 (22.3) | 73.5 (16.8) | 71.7 (16.1) | 80.8 |
| ErG | 167.4 (30.2) | 99.1 (22.7) | 67.4 (13.9) | 80.8 |
| RDkitDesc | 122.3 (37.2) | - | 73.9 (11.2) | 80.8 |
| MACCSKeys | 145.6 (25.0) | 81.1 (16.4) | 72.5 (11.8) | 80.8 |
| Estate | 165.7 (24.5) | 93.5 (22.3) | 78.5 (16.3) | 80.8 |
| $\Delta H_{\text{endo(max)}}$ (kJ/mol) | | | | |
| Mordred | $2.4\times10^7$ ($4.2\times10^7$) | - | 5.3 (1.3) | 4.6 |
| ECFP | 4.8 (1.7) | 4.8 (2.1) | 4.2 (1.4) | 4.6 |
| Avalon | 4.3 (2.0) | 4.8 (2.3) | 5.3 (1.4) | 4.6 |
| ErG | 5.8 (1.4) | 4.7 (2.1) | 4.5 (1.3) | 4.6 |
| RDkitDesc | 10.1 (9.5) | - | 5.2 (1.5) | 4.6 |
| MACCSKeys | 4.8 (1.2) | 4.7 (2.6) | 4.6 (1.4) | 4.6 |
| Estate | 5.9 (1.8) | 4.8 (1.8) | 4.4 (1.3) | 4.6 |
| $\Delta H_{\text{endo(min)}}$ (kJ/mol) | | | | |
| Mordred | $1.9\times10^7$ ($2.3\times10^7$) | - | 5.1 (1.3) | 3.7 |
| ECFP | 4.3 (1.5) | 3.9 (1.9) | 4.5 (1.8) | 3.7 |
| Avalon | 4.0 (1.3) | 3.9 (1.6) | 4.9 (1.4) | 3.7 |
| ErG | 5.6 (1.9) | 3.9 (2.3) | 3.8 (1.4) | 3.7 |
| RDkitDesc | 11.5 (11.1) | - | 4.1 (1.2) | 3.7 |
| MACCSKeys | 4.2 (2.0) | 3.9 (2.0) | 4.3 (1.7) | 3.7 |
| Estate | 4.6 (1.3) | 4.2 (1.8) | 4.0 (1.8) | 3.7 |
| $\Delta H_{\text{exo(max)}}$ (kJ/mol) | | | | |
| Mordred | $1.2\times10^7$ ($1.9\times10^7$) | - | 3.7 (1.0) | 3.2 |
| ECFP | 3.8 (0.8) | 3.5 (1.6) | 3.1 (1.2) | 3.2 |
| Avalon | 3.3 (1.1) | 3.5 (1.5) | 3.2 (1.1) | 3.2 |
| ErG | 6.5 (2.2) | 3.6 (1.6) | 3.5 (1.0) | 3.2 |
| RDkitDesc | 12.6 (11.5) | - | 3.9 (1.2) | 3.2 |
| MACCSKeys | 4.6 (1.3) | 3.5 (1.5) | 3.8 (1.0) | 3.2 |
| Estate | 5.0 (1.6) | 3.6 (1.4) | 3.7 (0.9) | 3.2 |

| $\Delta H_{exo(min)}$ (kJ/mol) | | | | |
|---|---|---|---|---|
| Mordred | $1.5 \times 10^7$ ($1.9 \times 10^7$) | - | 4.3 (1.1) | 3.9 |
| ECFP | 3.4 (0.9) | 4.1 (1.7) | 3.4 (1.3) | 3.9 |
| Avalon | 3.2 (1.0) | 4.1 (1.9) | 3.3 (1.2) | 3.9 |
| ErG | 6.4 (1.9) | 4.1 (1.8) | 3.9 (1.5) | 3.9 |
| RDkitDesc | 8.7 (3.6) | - | 4.3 (1.0) | 3.9 |
| MACCSKeys | 4.5 (1.1) | 4.1 (2.1) | 4.5 (1.3) | 3.9 |
| Estate | 5.1 (2.0) | 4.2 (1.2) | 4.5 (1.3) | 3.9 |

**Table S5. Top-10 ranked important features in the regression of each property.** The value in bracket is the averaged value of feature importance.

| Rank | $T_{endo(max)}$ | $T_{endo(min)}$ | $T_{exo(max)}$ | $T_{exo(min)}$ |
|---|---|---|---|---|
| 1 | **VSA_EState4** (0.044) | **ATSC2d** (0.044) | **ATSC1are** (0.022) | ATSC1se (0.039) |
| 2 | ATSC3p (0.034) | **AATS0p** (0.036) | **VSA_EState4** (0.021) | **VSA_EState4** (0.037) |
| 3 | **ATSC3v** (0.031) | **VSA_EState4** (0.034) | ATSC5dv (0.020) | **ATSC1are** (0.036) |
| 4 | SLogP_VSA2 (0.031) | JGI1 (0.030) | AATSC3m (0.019) | GATS1d (0.023) |
| 5 | **ATSC2d** (0.018) | AATSC1dv (0.028) | Xc-3d (0.018) | ATSC5d (0.017) |
| 6 | **AATS0p** (0.017) | GATS1p (0.023) | AATSC3Z (0.017) | JGI1 (0.017) |
| 7 | GATS3p (0.016) | **ATSC3v** (0.020) | VSA_Estate3 (0.014) | ATSC1pe (0.017) |
| 8 | GATS2d (0.015) | ATSC1se (0.019) | SssssC (0.014) | JGT10 (0.015) |
| 9 | Xc-3dv (0.013) | GATS3d (0.017) | Xc-3dv (0.013) | GATS1p (0.015) |
| 10 | GATS3v (0.013) | ATSC1pe (0.017) | EState_VSA1 (0.011) | AATSC0d (0.014) |

**Red bold**: the descriptor commonly appeared in $T_{endo(max)}$, $T_{endo(min)}$, $T_{exo(max)}$, and $T_{exo(min)}$

**Black bold**: the descriptor commonly appeared in $T_{endo(max)}$ and $T_{endo(min)}$

**Orange bold**: the descriptor commonly appeared in $T_{exo(max)}$ and $T_{exo(min)}$

**Supplementary references**

1. K. Moovendaran and S. Natarajan, *Spectrochim. Acta, Part A*, 2015 **135**, 317-320.

2. P. F. Li, Y. Y. Tang, Z. X. Wang, H. Y. Ye, Y. M. You and R. G. Xiong, *Nat. Commun.*, 2016, **7**, 13635.

3. A. E. Frumkin, N. V. Yudin, K. Y. Suponitsky and A. B. Sheremetev, *Mendeleev Commun.*, 2018, **28**, 135-137.

4. E. Nauha and J. Bernstein, *J. Pharma. Sci.*, 2015, **104**, 2056-2061.

5. A. Lemmerer, *CrystEngComm*, 2012, **14**, 2465-2478.

6. M. Rafilovich, J. Bernstein, M. B. Hickey and M. Tauber, *Cryst. Growth Des.*, 2007, **7**, 1777-1782.

7. C. S. Yang, Y. H. Tan, C. F. Wang, S. P. Chen, B. Wang, H. R. Wen and Y. Z. Tang, *Chem. Phys.*, 2018, **502**, 66-71.

8. T. Kusukawa, Y. Kojima, and F. Kannen, *Chem. Lett.*, 2019, **48**, 1213-1216.

9. C. Liu, J. Chen, C. Xu, H. Hao, B. Xu, D. Hu, G. Shi and Z. Chi, *Dyes Pigm.*, 2020, *174*, 108093.

10. A. Ainurofiq, R. Mauludin, D. Mudhakir, D. Umeda, S. N. Soewandhi, O. D. Putra and E. Yonemochi, *Eur. J. Pharma. Sci.*, 2018, **111**, 65-72.

11. R. Bhowal, A. A. Balaraman, M. Ghosh, S. Dutta, K. K. Dey and D. Chopra, *J. Am. Chem. Soc.*, 2020, **143**, 1024-1037.

12. V. Kumar, R. Thaimattam, S. Dutta, P. Munshi and A. Ramanan, *CrystEngComm*, 2017, **19**, 2914-2924.

13. T. Kojima, F. Kato, R. Teraoka, Y. Matsuda, S. Kitagawa and M. Tsuhako, *Chem. Pharma. Bull.*, 2007, **55**, 407-411.

14. S. Ying, M. Chen, Z. Liu, M. Zheng, H. Zhang, S. Xue and W. Yang, *J. Mater. Chem. C*, 2017, **5**, 5994-5998.

15. A. Kapf, H. Eslahi, M. Blanke, M. Saccone, M. Giese and M. Albrecht, *New J. Chem.*, 2019, **43**, 6361-6371.

16. P. Rani, A. Husain, A. Shukla, N. Singla, A. K. Srivastava, G. Kumar, K. K. Bhasin and G. Kumar, *CrystEngComm*, 2021, **23**, 1859-1869.