1	Supplementary Information for:
2	Active Learning-Guided Exploration of Thermally Conductive
3	Polymers Under Strain
4	Renzheng Zhang ¹ , Jiaxin Xu ¹ , Hanfeng Zhang ¹ , Guoyue Xu ¹ , Tengfei Luo ^{1,2*}
5	1. Department of Aerospace and Mechanical Engineering, University of Notre Dame, Notre
6	Dame, Indiana 46556, United States
7	2. Department of Chemical and Biomolecular Engineering, University of Notre Dame, Notre
8	Dame, Indiana 46556, United States
9	3. Lucy Family Institute for Data & Society, University of Notre Dame, Notre Dame, Indiana
10	46556, United States
11	

12 * Corresponding email: tluo@nd.edu

13

Supplementary Note 1. Optimization process in MD simulation

In the optimization process, the polymer systems were initially simulated with 14 15deactivated electrostatic interactions, while a 0.300 nm cutoff was implemented for Lennard-16 Jones interactions. The system was first subjected to a simulation under the NPT (constant 17number of atoms, pressure, and temperature) ensemble at 100 K for 2 ps, with a 0.1 fs timestep 18 applied. The system was then gradually heated from 100 K to 1000 K over 1 ns under the NVT 19 ensemble, followed by a simulation under the NPT ensemble for 50 ps at 0.1 atm and 1000 K. 20 This stage was intended to further relax the structure and ensure the complete eradication of 21 close contacts. Subsequently, the system underwent further simulation under NPT at 1000 K 22 for 1 ns, during which the pressure was permitted to rise from 0.1 atm to 500 atm, with a 1 fs 23 timestep and SHAKE constraints ¹ implemented. The SHAKE constraint facilitates the usage 24 of a larger 1 fs timestep in the presence of lightweight hydrogen atoms within the system. Throughout these stages, the deactivation of electrostatic interactions and the use of a range of 2526 ad-hoc simulation procedures were adopted to eliminate close atomic contacts, a vital measure 27 to prevent system disruption in subsequent simulation steps.

These simulation procedures are collectively referred to as the initialization process. The obtained polymer system was then annealed with activated electrostatic interactions, using the PPPM (Particle–Particle–Particle–Mesh)-based Ewald sum method for computation. A 0.800 nm cutoff was used for the Lennard-Jones interactions. In this annealing phase, the system was first run in an NPT ensemble at 1 atm and 1000 K for 2 ps using a 0.1 fs timestep. Subsequently, it was cooled from 1000 K to 300 K at a rate of 140 K/ns, under an NPT ensemble at 1 atm. This was followed by another NPT run at 300 K and 1 atm for 8 ns, utilizing
 a 1 fs timestep and SHAKE constraints to achieve the final amorphous state. Such a procedure

has been shown to yield converged thermal conductivity (TC) results with NEMD².

37

36

38 Supplementary Note 2. Gaussian process regression model

39 A GPR model is a nonparametric Bayesian method that models the distribution over functions. 40 In a GPR framework, we assume the observed data y consists of a latent function $f(\mathbf{x})$ corrupted by independent and identically normally distributed noise ϵ , where $\epsilon \sim$ 41 42 $N(0,\sigma^2)$. The function f is modeled as a Gaussian process, meaning that any finite set of function values follows a multivariate normal distribution with a specified mean function m(x)43 and a kernel (covariance) function k(x, x'). The kernel encodes the similarity between 44 different input points. Given a dataset D, the joint distribution over the observed outputs y and 45 46 the function values at a new point x_* is:

47
$$\begin{bmatrix} y \\ f(\mathbf{x}_*) \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} m(\mathbf{X}) \\ m(\mathbf{x}_*) \end{bmatrix}, \begin{bmatrix} K(\mathbf{X}, \mathbf{X}) + \sigma^2 I & k(\mathbf{X}, \mathbf{x}_*) \\ k(\mathbf{x}_*, \mathbf{X}) & k(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix}\right)$$

Here, $K(\mathbf{X}, \mathbf{X})$ is the covariance matrix between the training points, and $k(\mathbf{X}, \mathbf{x}_*)$ is the covariance vector between the training points and the new point. Using Bayes' rule, we can derive the posterior distribution for the predicted mean $\mu_*(\mathbf{x}_*)$ and the variance $\sigma^2_*(\mathbf{x}_*)$, given by:

52
$$\mu_*(\mathbf{x}_*) = m(\mathbf{x}_*) + k(\mathbf{x}_*, \mathbf{X})^T [K(\mathbf{X}, \mathbf{X}) + \sigma^2 I]^{-1} (y - m(\mathbf{X}))$$

53
$$\sigma_*^2(\mathbf{x}_*) = k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, \mathbf{X})^T [K(\mathbf{X}, \mathbf{X}) + \sigma^2 I]^{-1} k(\mathbf{X}, \mathbf{x}_*)$$

For the kernel function k(x, x'), we use a Rational Quadratic kernel combined with a White kernel to account for noise in the observations. The Rational Quadratic Kernel is chosen because of it can model both small- and large-scale variations in the data, and the White Kernel is chosen to account for observational noise, ensuring that the model is robust to noisy data and does not overfit small fluctuations due to measurement errors. The combined kernel function is defined as:

60
$$k(\mathbf{x}, \mathbf{x}') = (1 + \frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\alpha \cdot l^2})^{-\alpha} + \sigma^2 \cdot \delta(\mathbf{x} - \mathbf{x}')$$

61 where, $||\mathbf{x} - \mathbf{x}'||$ is the squared Euclidean distance between input points \mathbf{x} and \mathbf{x}' , α controls 62 the relative weighting of large and small-scale variations in the data, l is the characteristic 63 length scale, determining how quickly the function varies with the input. σ^2 represents the 64 noise variance, and $\delta(\mathbf{x} - \mathbf{x}')$ is the Kronecker delta function, which ensures that noise only 65 affects the diagonal of the covariance matrix.

66

67 Supplementary Note 3. Bayesian optimization process

Fig. S1 shows the parity plots of each iteration and GPR prediction & uncertainty on the PoLyInfo database after iteration 2. For each iteration, the newly incorporated data are marked by red circles. The predictive accuracy, estimated by the R², is above 0.8 in all iterations. The upper-bound GPR-predicted TC increases overall and reaches 1.1358 W/m·K in iteration 8. The amount of polymers that have GPR-predicted TC over 1.0 W/m·K increases from 0 in iteration 2 to 485 in iteration 8. **Table S1** provides detailed information, including PID, MDlabeled TC and SMILES, of the 35 BO-recommended polymers.



Fig. S1 | GPR model performances from iteration 3 to iteration 8. Parity plot between GPR-76 predicted TC versus MD-labeled TC and GPR-predicted TC mean and uncertainty of each 77polymer in the PoLyInfo database in a iteration 3, b iteration 4, c iteration 5, d iteration 6, e 78 iteration 7, **f** iteration 8, with the five new BO-guided data marked by red circles. 7980

Iteration	PID	TC (W/m·K)	Smiles
	P432586	0.348	*c1ccc(Oc2ccc(-c3nc4ccc(Oc5ccc(N6C(=O)c7ccc(C(c8ccc9c(c8)C(=O)N(c8ccc(Oc%10ccc%11nc(*)c(- c%12ccccc%12)nc%11c%10)cc8)C9=O(C(F)(F)F)C(F)(F)F)cc7C6=O(cc5)cc4nc3-c3ccccc3)cc2)cc1
	P460108	0.2	*/C(=N\c1ccccc1)N(*)c1ccccc1
Random	P010100	0.223	*C1CCC1*
	P040310	0.238	*CC(*)(CC(=0)OCCCCO)C(=0)OCCCCCC
	P100773	0.521	*CCCCNC(=0)CCCCCCC(=0)NCCCCNC(=0)C(=0)N*
	P120025	0.801	*CCCCCCCCCCCCCC(=0)NCCCCCCCCCCC(=0)N*
	P400028	0.726	*CCCCCCCCCCCCC(=0)CCCCCCCC(=0)NCCCCCCCCCCCC(=0)C(=0)N*
Iter_2	P340177	0.525	*CC(*)(C)C(=0)NCCCCCCCC(=0)NCCCCCCCC(=0)OC
	P340179	0.527	*CC(*)(C)C(=O)NCCCCCCCCC(=O)NCCCCCCCC(=O)O
	P100719	0.605	*CCCCCCNC(=0)C(CCCCCCCCCCCCCCCCC)C(=0)N*
	P100029	0.926	*CCCCCCCCCCCCCCC(=0)CCCCCCCCCCCCC(=0)N*
	P100235	0.741	*CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
Iter_3	P100238	0.726	*CCCCCCCCCCCCCC(=0)CCCCCCCCCCCCCCCCC(=0)N*
	P332306	0.743	*0000000000000000(*)0000000000000000000
	P402128	1.086	*CCCCCCCCCCCCC(=0)CCCCCCCCCCCCCCC(=0)N*
	P090293	0.517	*CCCCCC(=0)NCCCCCCNC(=0)CCCCCCC(=0)O*
	P090368	0.454	*CCCCCCCCC(=0)CCCCCNC(=0)CCCCCC(=0)NCCCCCC(=0)O*
Iter_4	P100747	0.709	*CCCCCCCCCC(=0)CCCCCCCC(=0)NCCCCCCCCCC(=0)C(=0)N*
	P100809	0.651	*CCCCCCCCC(=0)CCCCCCCC(=0)NCCCCCCCCC(=0)C(=0)N*
	P400027	0.693	*CCCCCCCCC(=0)CCCCCCCC(=0)NCCCCCCCCC(=0)C(=0)N*
	P090253	0.888	*0(0=)20222220(0=)20222222222222222222222
	P100030	1.136	*CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
Iter_5	P100696	0.919	*CCCCCCCCCCCCCCCCCC(=0)N*
	P332076	0.741	*/C=C(/*)CCCCCCCCCCCCCCCC(=0)0
	P372661	0.593	*CCOCCOCCOCc1cc(CO*)cc(OCCCCCCCCCCCCCCC)c1
	P332081	0.728	*/C=C(/*)C#CCCCCCCCCCCCCCCCC(=O)O
	P080116	0.689	*CCCCCCCC(=O)NCCc1ccc(CCNC(=O)CCCCCCCS*)cc1
Iter_6	P100124	0.717	*CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
	P100063	0.658	*CCclccc(CCNC(=O)CCCCCCCCCCCCCC(=O)N*)ccl
	P390132	0.948	*CCCCCCCCCCCCCCC(=0)CCCCCCC(=0)O*
	P372375	0.429	*UUUNU(=U)UNU(=U)U(NU(=U)U(Cc1ccccc1)NC(=U)UUUUU(=U)NU(Cc1ccccc1)C(=U)NC(C(= 0)NCC(=O)NCCCOCCOCCO*)C(C)CC(C)C
	P402222	1.036	*CCCCCCCCCCCCC(=0)CCCCCCCCCC(=0)N*
Iter_7	P402250	1.09	*CCCCCCCCCCCC(=0)CCCCCCCCCCCCCCC(=0)N*
	P110210	0.839	*CCCCCCCCCCCCCC(=O)CCCCCCCCCC(=O)N*
	P392523	0.824	*0(0=))))(0=))))))))))))))))))))))))))))
	P450068	0.263	*CC(OC(=0)OC1CCC2(C)C(=CCC3C2CCC2(C)C(C(C)CCCC(C)C)CCC32)C1)C(COC(=0)O*)OC(=0)O C1CCC2(C)C(=CCC3C2CCC2(C)C(C)C(C)CCCC32)C1
	P402218	0.884	*CCCCCCCCC(=0)CCCCCCCC(=0)CCCC(=0)N*
Iter 8	P100754	0.878	*CCCCCNC(=0)CCCCCCCCCCCCCC(=0)N*
	P522013	0.135	*CC(*)(F)C(=O)OCC(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)(F)(F)(F)(F)(F)(F)(F)(F)(F)(F)(F)(
	P382405	0.468	*clnc(CCCCCCC)c(-c2sc(-c3sc(-c4cc(CCCCCCC)c(*)s4)cc3CCCCCCCC)nc2CCCCCC)s1

Table S1. MD-labeled TC of 35 BO-guided polymers and their SMILES

81

82 Supplementary Note 4. Comparative analysis of other ML models

To compare the performance of the GPR model, we test two other popular ML models: random 83 84 forest (RF) and gradient boosting regressor (GBR), both of which are insensitive to dimensionality, so no dimensionality reduction process is performed on both models. Fig. 85 86 S2a&d show the parity plot between RF-predicted TC and MD-labeled TC in iteration 1 and 87 iteration 8. Fig. S2b&e show the parity plot between GPR-predicted TC and RF-predicted TC in iteration 1 and iteration 8, where R^2 are 0.72 and 0.53, respectively. Since neither the RF 88 89 model nor the GBR model provides the uncertainty, we only have the predicted value of TC, 90 as shown in Fig. S2c&f. Similarly, Fig. S2g&j show the parity plot between GBR-predicted 91 TC and MD-labeled TC in iteration 1 and iteration 8. Fig. S2h&k show the parity plot between 92 GPR-predicted TC and GBR-predicted TC in iteration 1 and iteration 8, where R² are 0.64 and 93 0.52, respectively. The predicted value of TC in iteration 1 and iteration 8 are shown in Fig. 94 S2i&I. It is noteworthy that with the initial 36 training data, the distributions of the RF and 95 GBR models in iteration 1 are more similar to the GPR model, where the agreement is more reasonable than that in iteration 8. With more data incorporated, the R² of RF-GPR and GBR-96 97 GPR decreased instead. It has been validated that the GPR model has accurate predicted TC 98 values that match the MD simulation, especially for high TC polymers. However, the 99 boundaries of RF and GBR models are both around 1, which means both RF and GBR models 100 underestimate the strained polymers TC to some extent (more datapoints are above the gray 101 dashed line).



103 Fig. S2 | Comparative analysis of RF and GBR models. a Parity plot between RF-predicted TC and MD-labeled TC in iteration 1. b Parity plot between GPR-predicted TC and RF-104 105predicted TC in iteration 1. c RF-predicted TC of each polymer in the PoLyInfo database in 106 iteration 1. d Parity plot between RF-predicted TC and MD-labeled TC in iteration 8. e Parity 107 plot between GPR-predicted TC and RF-predicted TC in iteration 8. f RF-predicted TC of each 108 polymer in the PoLyInfo database in iteration 8. g Parity plot between GBR-predicted TC and 109 MD-labeled TC in iteration 1. h Parity plot between GPR-predicted TC and GBR-predicted TC in iteration 1. i GBR-predicted TC of each polymer in the PoLyInfo database in iteration 1. 110

102

- 111 j Parity plot between GBR-predicted TC and MD-labeled TC in iteration 8. k Parity plot
- between GPR-predicted TC and GBR-predicted TC in iteration 8. I GBR-predicted TC of each
- 113 polymer in the PoLyInfo database in iteration 8.







Fig. S3 | Summary of the 30 MD-labeled high TC (>0.8 W/m·K) polymers' categories.

124 Supplementary Note 6. Herman's orientation factor

Herman's orientation factor $(f)^{3,4}$ is used to characterize the segment alignment along different 125directions: $f = 1.5 < cos^2 \theta > -0.5$, where θ is the angle between a carbon-carbon bond and 126 the reference direction. A value of f = -0.5 means that the orientation is perpendicular to the 127 128 selected direction, and f = 1 means that the orientation is parallel to the selected direction. A 129 value of 0 indicates a completely random orientation. We select polymer P110210 to simulate its f and TC in the oriented and perpendicular directions at different draw ratio. Fig. S4a shows 130 f increases with larger draw ratio along the draw direction, while decreasing in the 131 132perpendicular direction. Fig. S4b, respectively, indicates the changes of TC in the oriented and 133perpendicular direction at different draw ratio, which is consistent with Fig. S4a. Fig. S4c visualizes the structures of this polymer under different strain ratios. The main chain carbon 134135atoms, indicated in yellow, are more aligned with the draw ratio increasing.



Fig. S4 | Herman's orientation factor simulation for polymer P110210 at different draw
ratio. a Averaged Herman's orientation factor and b TC in oriented and perpendicular
directions at different draw ratios of polymer P110210. c Structures visualization of polymer

140 P110210 at different draw ratios.

141 Supplementary Note 7. Structural comparison of high TC and low TC polymers

- Fig. S5 shows the polymers that have linear chain structures (right) and non-linear chain structures (left). For linear chain polymers, they have much higher TC and f than non-linear
- 144 chain polymers due to the higher orientation factor and chain alignments.





147 **Reference**

- 148 (1) Ryckaert JP, Ciccotti G, Berendsen HJC. Numerical integration of the cartesian equations
- 149 of motion of a system with constraints: molecular dynamics of n-alkanes. J Comput Phys.
- 150 1977; 23(3):327–341.
- 151 (2) Ma R, Zhang H, Xu J, Sun L, Hayashi Y, Yoshida R, Shiomi J, Wang J, Luo T. Machine
- 152 learning-assisted exploration of thermally conductive polymers based on high-throughput
- 153 molecular dynamics simulations. Mater Today Phys. 2022;28:100850.
- 154 (3) Zhang T, Luo T. Role of chain morphology and stiffness in thermal conductivity of
- 155 amorphous polymers. J Phys Chem B. 2016;120(4):803–812.
- 156 (4) Pal S, Balasubramanian G, Puri IK. Modifying thermal transport in electrically conducting
- polymers: effects of stretching and combining polymer chains. J Chem Phys.
 2012;136(4):044901–044907.

159