

1

## 2 **Supplementary Information for**

### 3 **Archetypal Landscapes for Deep Neural Networks**

4 **P.C. Verpoort, A.A. Lee, D.J. Wales**

5 **Philipp C. Verpoort.**

6 **E-mail: [pcv22@cam.ac.uk](mailto:pcv22@cam.ac.uk)**

#### 7 **This PDF file includes:**

- 8     Supplementary text
- 9     Figs. S1 to S5
- 10    Tables S1 to S2
- 11    Caption for Movie S1
- 12    References for SI reference citations

#### 13 **Other supplementary materials for this manuscript include the following:**

- 14     Movie S1

## 15 Supporting Information Text

16 In this supporting information, we explain the data set used for training and testing in the main contri-  
17 bution. Sec. S1 outlines the LJAT<sub>3</sub> classification problem that was employed and the data sets that were  
18 generated. Sec. S2 describes the visualisation of solutions that is used throughout our main report. We also  
19 tabulate mean values for uphill and downhill barriers for individual transition states and for the pathway  
20 with the lowest maximum transition state energy connecting other minima to the global minimum.

### 21 S1. Predicting the Outcome of Geometry Optimisation for an Atomic Cluster

22 This benchmarking problem has been used in several previous contributions that employed neural network  
23 fits with single hidden layers (1–3). This work investigated how the corresponding machine learning  
24 landscapes and predictions varied with the number of nodes and the number of training data, including  
25 the effect of memory in sequences of molecular configurations. The system is a triatomic cluster bound  
26 by pairwise Lennard-Jones (4) and three-body Axilrod–Teller (5) terms, parameterised so that there are  
27 three permutational isomers of a linear minimum, distinguished by the central atom, and one additional  
28 minimum for an equilateral triangle with  $D_{3h}$  symmetry. The total potential energy for this LJAT<sub>3</sub> cluster  
29 is

$$30 \quad V = 4\epsilon \sum_{i<j} \left[ \left( \frac{\sigma}{r_{ij}} \right)^{12} - \left( \frac{\sigma}{r_{ij}} \right)^6 \right] + Z \sum_{i<j<k} \left[ \frac{1 + 3 \cos \theta_1 \cos \theta_2 \cos \theta_3}{(r_{ij}r_{ik}r_{jk})^3} \right], \quad [1]$$

31 where  $\theta_1$ ,  $\theta_2$  and  $\theta_3$  are the internal angles of the triangle formed by atoms  $i$ ,  $j$ ,  $k$ .  $r_{ij}$  is the distance  
32 between atoms  $i$  and  $j$ , and  $Z$  is a parameter that weights the contribution of the three-body term. For  
33  $Z = 2$  the linear minima have potential energy  $V = -2.219\epsilon$ , and the triangle lies slightly higher with  
34  $V = -2.185\epsilon$ . For the triangle  $r_{12} = r_{13} = r_{23} = 1.16875\sigma$ , and in the linear minima the nearest-neighbour  
35 distances are both  $1.10876\sigma$ .

36 The aim of this multinomial logistic regression problem is to predict which of the four local minima  
37 a geometry optimisation will find, given some information about initial or intermediate configurations  
38 in terms of the interparticle distances. To generate data we consider starting geometries constructed  
39 from randomly distributing the three atoms in a cube of side length  $L$ . The initial values of  $r_{12}$  and  $r_{13}$   
40 were employed as the input data for all the tests conducted in the present work, and we considered two  
41 datasets, the first (D1) for 10,000 minimisations for a cube with  $L = 2\sqrt{3}\sigma$ , and the second (D2) for  
42 200,000 minimisations with  $L = 1.385\sigma$ . We used half of each dataset for training and half for testing,  
43 where appropriate. Dataset D1 was employed in previous work, while D2 was generated for the present  
44 investigation. Each local minimisation employed the same LBFGS minimisation routine described in  
45 **Methods**, and the convergence condition was taken as  $10^{-6}$  for the root mean square gradient in reduced  
46 units of  $\epsilon/\sigma$ .

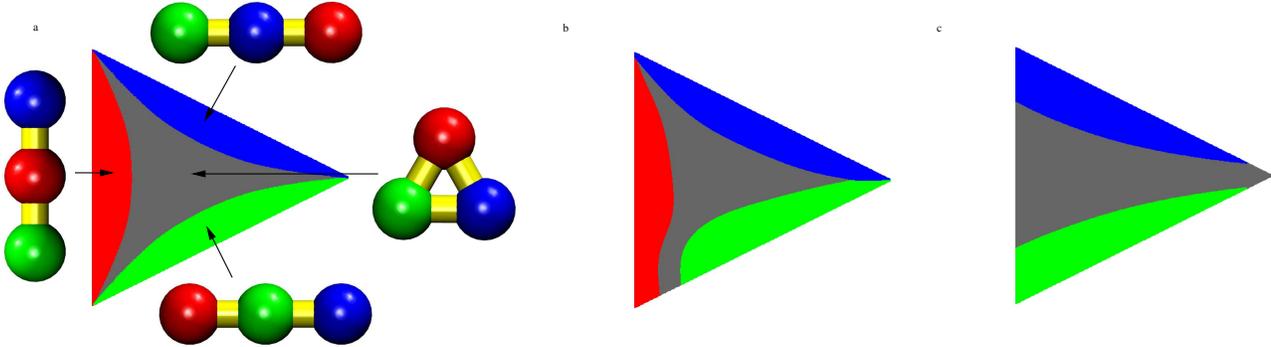
47 The molecular configuration is completely characterised by three interparticle distances,  $r_{12}$ ,  $r_{13}$ , and  
48  $r_{23}$ . If we supply sufficient training data with these three inputs, predicting the outcome of minimisation  
49 can be essentially perfect. The problem is then equivalent to learning the basin of attraction for each local  
50 minimum, which is a well-defined volume of configuration space for steepest-descent minimisation (6, 7).

51 By restricting the input data to  $r_{12}$  and  $r_{13}$ , and omitting  $r_{23}$ , we make the prediction problem harder.  
52 For the linear minima with atom 2 or atom 3 in the middle,  $r_{13}$  and  $r_{12}$ , respectively, are much larger  
53 than for the triangle. However, these distances are only about 5% different in the triangle and the linear  
54 minimum with atom 1 in the central position. The basins of attraction of the triangle and this third  
55 linear minimum therefore overlap significantly in the space defined by  $r_{12}$  and  $r_{13}$ . The best predictions  
56 we can achieve will therefore occur when we have converged the relative probabilities of finding these two  
57 structures as a function of  $r_{12}$  and  $r_{13}$ .

58 The same considerations will apply for larger molecules: if we sample the whole configuration space  
59 sufficiently, we should be able to predict which basin of attraction any starting structure corresponds to,

60 and the corresponding local minimum. Otherwise, the best we can do is to learn the relative probabilities,  
 61 averaged over the missing degrees of freedom. Hence we find the current benchmark appealing because  
 62 of the ability to generate arbitrary amounts of training and testing data, because of the clear physical  
 63 interpretations, and because of the practical importance of the configuration volumes themselves. For  
 64 example, the volume of basins of attraction provides measures of configurational entropy, which has been  
 65 applied to analyse granular packings (8). Reliable prediction for the outcome of minimisation would enable  
 66 us to reduce the time required for such calculations by stopping earlier, with a weaker convergence threshold  
 67 on the magnitude of the gradient (9).

## 68 S2. Visualisation of Solutions



**Fig. S1.** Graphical representation of the LJAT<sub>3</sub> classification problem. (a) Colored according to the true outcome determined by geometry optimization for the LJAT<sub>3</sub> cluster. The four optimal atomic configurations are associated with their corresponding basins of attraction. (b) Colored according to the predictions for the global minimum of a single hidden layer neural network with 3 hidden nodes and 100 training data confined in the plane  $\mathcal{R}'$  (AUC 0.98 from corresponding test set). (c) Colored according to the predictions for the global minimum of a single hidden layer neural network with 10 hidden nodes trained on 100,000 training data in  $\mathcal{R}$  (AUC 0.79 from corresponding test set).

69 Some insight into the different local minima in the cost function for a given neural network and training  
 70 data can be obtained graphically for the LJAT<sub>3</sub> prediction problem (3, 10). We construct a two-dimensional  
 71 projection of coordinates in the plane  $r_{12} + r_{13} + r_{23} = 3r_e$  from the three-dimensional space  $\{r_{12}, r_{13}, r_{23}\}$ ,  
 72 where  $r_e = 2^{1/6}$  is the equilibrium bond length in a dimer and in the equilateral triangle minimum. The  
 73 orthogonal unit vectors  $\hat{\mathbf{v}}_1 = (1, 1, -2)/\sqrt{6}$  and  $\hat{\mathbf{v}}_2 = (1, -1, 0)/\sqrt{2}$  lie in this plane and are perpendicular to  
 74 the  $\{1, 1, 1\}$  direction. We define projected coordinates  $x = (r_{12} + r_{13} - 2r_{23})/\sqrt{6}$  and  $y = (r_{12} - r_{13})/\sqrt{2}$ . For  
 75 a regular  $700 \times 700$  grid with  $-\sqrt{3}r_e < x, y < \sqrt{3}r_e$ , we solve for  $r_{12}$ ,  $r_{13}$ , and  $r_{23}$  with  $r_{12} + r_{13} + r_{23} = 3r_e$ ,  
 76 which gives 79524 geometrically feasible  $(x, y)$  points, and the associated values of  $\{r_{12}, r_{13}, r_{23}\}$ , and  
 77 Cartesian coordinates. The 79524 values of  $r_{12}$  and  $r_{13}$  constitute a third dataset D3. The feasible  
 78 geometries are distributed over a triangle in  $(x, y)$  space, where the centre of each edge corresponds to a  
 79 linear geometry with two distances of  $3r_e/4$  and one of  $3r_e/2$ , and each vertex corresponds to two atoms  
 80 coincident and the third at a distance of  $3r_e/2$ . The equilateral triangular minimum maps to  $(x, y) = (0, 0)$ .

81 The pixels on the  $(x, y)$  grid are coloured according to the minimum with highest predicted probability  
 82 when the associated configuration is used as input data for any given neural network. If the equilateral  
 83 triangle has the highest probability the pixel is gray, while the three linear minima with atoms 1, 2, and  
 84 3 in the centre are coloured red, green and blue, respectively. Since we are omitting  $r_{23}$  from all the  
 85 training data, we anticipate that predictions will be significantly perturbed from previous calculations that  
 86 included all three distances as inputs (3). The target result is given by the known outcomes obtained by  
 87 energy minimisation with the same colour scheme (Fig. S1a). The basins of attraction for the three linear  
 88 minima are symmetrically disposed along the three edges of the triangle in the  $(x, y)$  projection, while the  
 89 remaining basin for the  $D_{3h}$  minimum has three-fold symmetry in this space.

90 Although this graphical representation only includes a subset of configurations in one plane defined in  
 91 the three-dimensional space  $\{r_{12}, r_{13}, r_{23}\}$ , comparison with the target reference pattern provides a very  
 92 useful indication of how well any particular neural network performs. It can be used for any set of weights

93 in any of the neural networks we consider, including transition states and all the configurations along the  
 94 pathways that connect the local minima. An illustrated profile for pathways leading to the global minimum  
 95 is shown in Fig. S3, and a movie with frames constructed from all the pathway configurations is available  
 96 as Supplementary Information. The capability to visualise cuts through a testing data set in terms of the  
 97 evolution in the predictive capabilities might prove useful in understanding how to construct better fits in  
 98 future work.

### 99 S3. Area Under Curve

100 To quantify the prediction capabilities of any given local minimum we calculated the area under curve  
 101 (AUC) for receiver operating characteristic (ROC) plots of the true positive rate,  $T_{\text{pr}}$ , against the false  
 102 positive rate,  $F_{\text{pr}}$ , as a function of the threshold probability,  $P$ , for predicting convergence to the equilateral  
 103 triangle. These rates are defined as

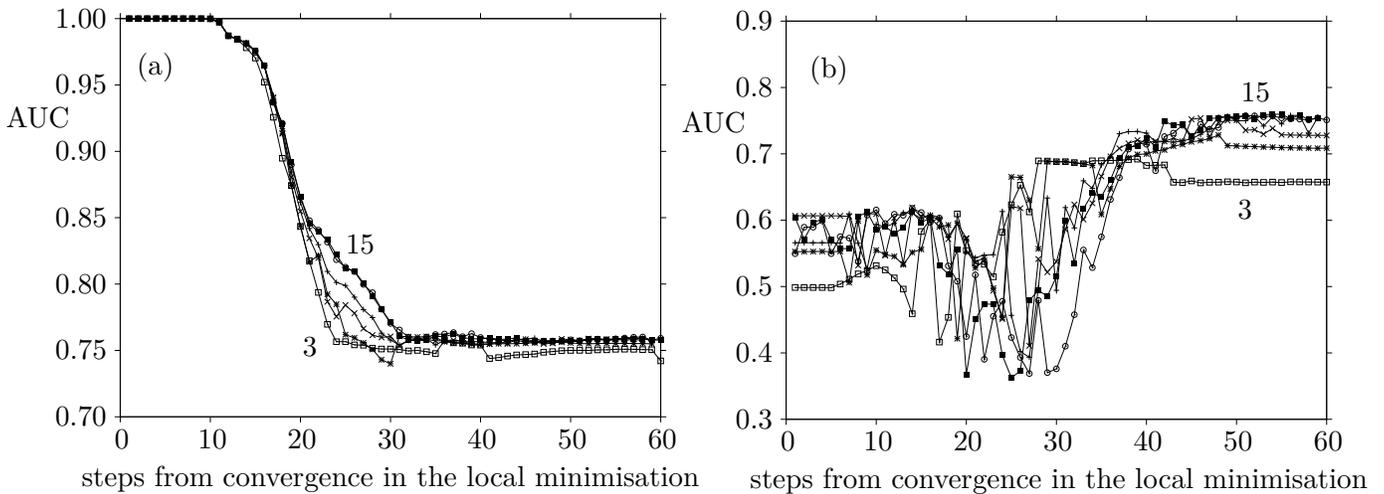
$$104 \quad T_{\text{pr}}(\mathbf{W}; P) = \frac{\sum_{d=1}^{N_{\text{data}}} \delta_{c(d),0} \Theta(p_0(\mathbf{W}) - P)}{\sum_{d=1}^{N_{\text{data}}} \delta_{c(d),0}},$$

$$105 \quad F_{\text{pr}}(\mathbf{W}; P) = \frac{\sum_{d=1}^{N_{\text{data}}} (1 - \delta_{c(d),0}) \Theta(p_0(\mathbf{W}) - P)}{\sum_{d=1}^{N_{\text{data}}} (1 - \delta_{c(d),0})}, \quad [2]$$

106 where  $\Theta$  is the Heaviside step function and  $\delta$  is the Kronecker delta. The AUC value is then

$$107 \quad \text{AUC}(\mathbf{W}) = \int_0^1 T_{\text{pr}}(\mathbf{W}; P) dF_{\text{pr}}(\mathbf{W}; P), \quad [3]$$

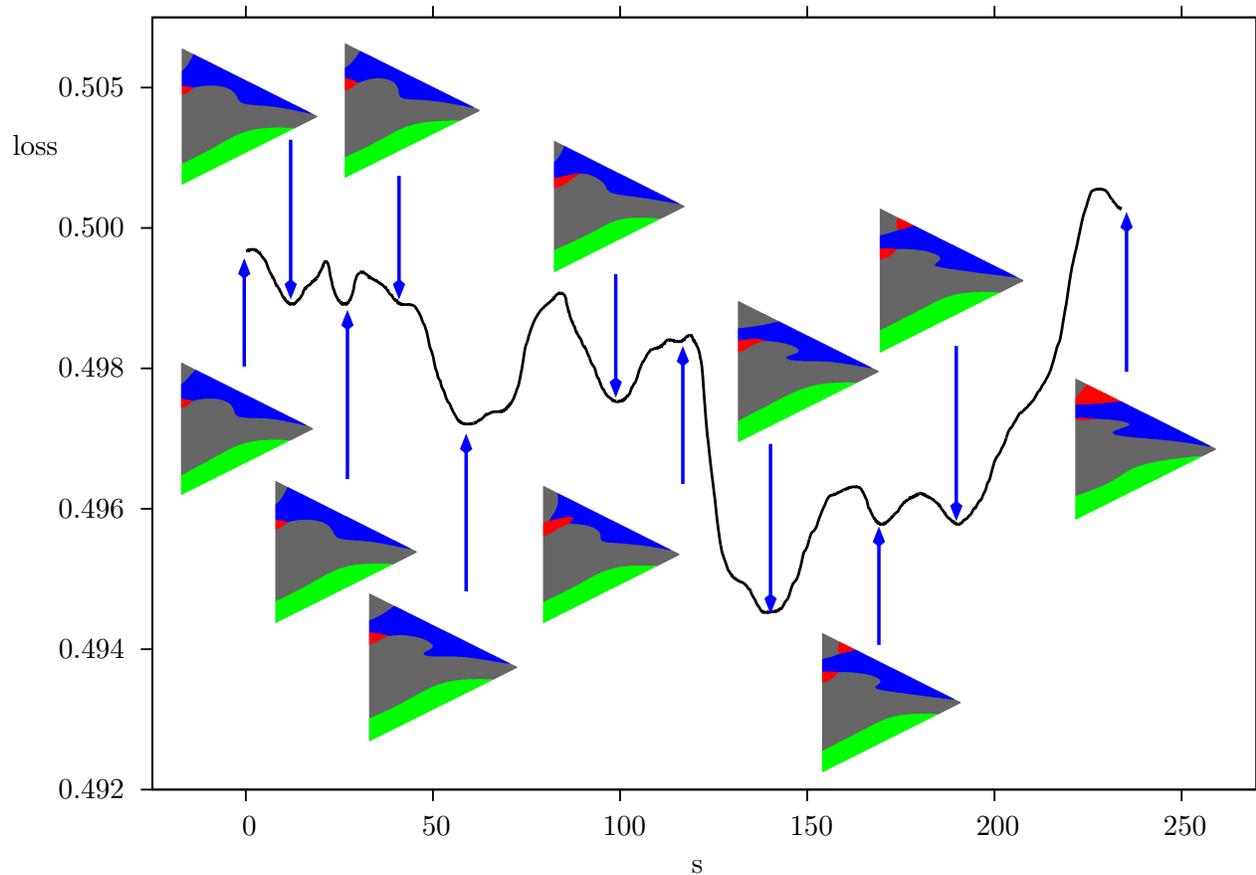
108 and was obtained numerically.



**Fig. S2.** AUC values obtained for the global minimum fit to 5000 training data in database D1 with networks containing a single hidden layer and 3, 4, 5, 6, 10 and 15 nodes. Global optimisation was performed for 5000  $(r_{12}, r_{13})$  pairs as a function of the number of steps to convergence in the geometry optimisation, as in previous work (2). (a) Results for the 5000 testing data in database D1 for  $(r_{12}, r_{13})$  pairs at the same number of steps to convergence as in each fit. (b) Results for the 79524 test data in database D3. The plots for 3 and 15 hidden nodes are indicated in each case, and the AUC values generally increase with the number of nodes as the training configurations approach the random initial configurations at the maximum number of steps from convergence.

109 We performed additional global optimisation runs for single hidden layers with 3, 4, 5, 6, 10 and 15 nodes  
 110 using configurations corresponding to saved  $(r_{12}, r_{13})$  data along the 5000 training minimisation sequences  
 111 in database D1. As in previous work (2), we find that the AUC values for the minima obtained in training,  
 112 and for configurations in the 5000 testing sequences in D1 at the same position in the minimisation, improve  
 113 systematically as the geometry optimisations approach convergence (Fig. S2a). However, if we apply the  
 114 solutions to the testing data in database D3 the best AUC values of around 0.75 correspond to fits using

115 the starting configurations, i.e. the random  $(r_{12}, r_{13})$  training data (Fig. S2b). Not surprisingly, training  
116 only on configurations close to the four equilibrium geometries of the cluster produces best fits that do not  
117 generalise as well to different configurations.



**Fig. S3.** Loss profile for a pathway involving eleven minima and ten transition states for the 1HL network trained on 250 data for LJAT<sub>3</sub> geometry optimisations. A graphical representation of the predictions for the D3 test set is indicated for each minimum. The global minimum is the eighth in the sequence. The horizontal axis corresponds to the integrated path length,  $s$ , which is calculated by treating the variable weights in the neural network according to a Euclidean distance metric.

**Table S1. Average uphill and downhill barriers for all the transition states and directly connected minima located in training, excluding degenerate rearrangements (7, 11). The second table is the mean of the barrier divided by the loss difference between the two minima, yielding a dimensionless parameter.**

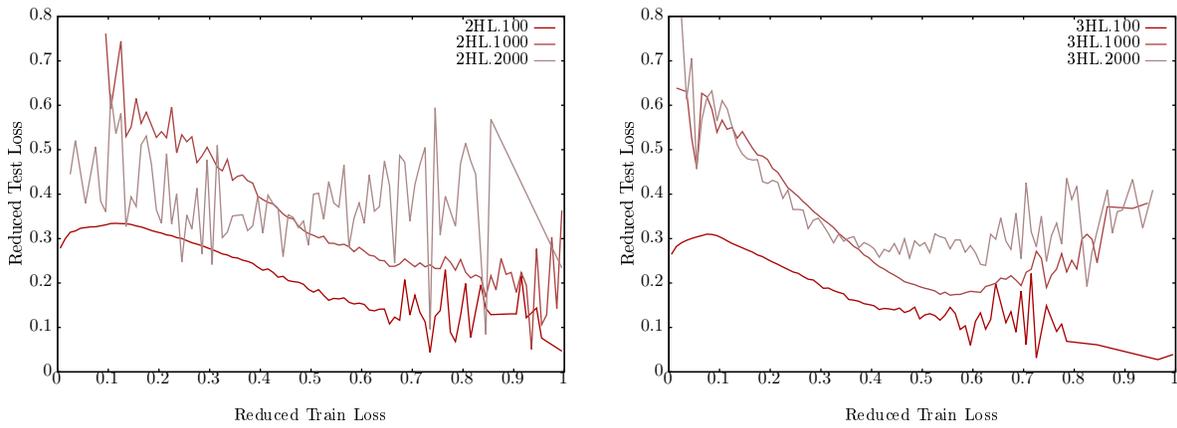
$N_{\text{data}}$	mean barriers					
	$H = 1$		$H = 2$		$H = 3$	
	uphill	downhill	uphill	downhill	uphill	downhill
100	$0.114 \times 10^{-1}$	$0.531 \times 10^{-2}$	$0.155 \times 10^{-1}$	$0.468 \times 10^{-2}$	$0.158 \times 10^{-1}$	$0.489 \times 10^{-2}$
1000	$0.985 \times 10^{-4}$	$0.196 \times 10^{-4}$	$0.239 \times 10^{-2}$	$0.853 \times 10^{-3}$	$0.314 \times 10^{-2}$	$0.978 \times 10^{-3}$
2000	$0.704 \times 10^{-4}$	$0.415 \times 10^{-4}$	$0.126 \times 10^{-2}$	$0.534 \times 10^{-3}$	$0.144 \times 10^{-2}$	$0.489 \times 10^{-3}$
10000	$0.109 \times 10^{-3}$	$0.455 \times 10^{-4}$	$0.466 \times 10^{-3}$	$0.186 \times 10^{-3}$	$0.748 \times 10^{-3}$	$0.274 \times 10^{-3}$
100000	$0.603 \times 10^{-4}$	$0.432 \times 10^{-4}$	$0.185 \times 10^{-3}$	$0.973 \times 10^{-4}$	$0.775 \times 10^{-3}$	$0.198 \times 10^{-3}$

$N_{\text{data}}$	mean barriers divided by loss difference of minima					
	$H = 1$		$H = 2$		$H = 3$	
	uphill	downhill	uphill	downhill	uphill	downhill
100	10.70	9.70	6.683	5.683	7.612	6.612
1000	2.131	1.131	5.767	4.767	4.915	3.915
2000	2.743	1.743	6.570	5.570	4.466	3.466
10000	3.167	2.167	7.952	6.952	5.672	4.672
100000	13.359	12.359	6.588	5.588	2.096	1.096

**Table S2. Average downhill barrier to the global minimum for all the other minima located in training, and average of the downhill barrier divided by the loss difference between the two minima (scaled column).**

$N_{\text{data}}$	$H = 1$		$H = 2$		$H = 3$	
		scaled		scaled		scaled
100	$0.105 \times 10^{-3}$	0.0178	$0.276 \times 10^{-2}$	0.0783	$0.324 \times 10^{-2}$	0.0781
1000	$0.298 \times 10^{-5}$	0.0905	$0.357 \times 10^{-3}$	0.0562	$0.540 \times 10^{-3}$	0.0462
2000	$0.589 \times 10^{-4}$	0.9777	$0.101 \times 10^{-3}$	0.0676	$0.566 \times 10^{-4}$	0.0294
10000	$0.415 \times 10^{-5}$	0.0597	$0.316 \times 10^{-4}$	0.4342	$0.324 \times 10^{-4}$	0.0477
100000	$0.663 \times 10^{-5}$	0.3301	$0.332 \times 10^{-4}$	0.3019	$0.216 \times 10^{-4}$	0.0286



**Fig. S4.** Reduced test loss plotted against reduced train loss of minima of the LJAT loss function landscapes for  $H = 2, 3$  and  $N_{\text{data}} = 100, 1000, 2000$ . The train loss is divided into 100 intervals and the test loss is averaged over all minima found with train loss in the interval. The reduced train (test) loss is defined as  $L_{\text{red}}(L) = \frac{L - L_{\text{min}}}{L_{\text{max}} - L_{\text{min}}}$ , where  $L_{\text{max}}$  is the maximal and  $L_{\text{min}}$  is the minimal train (test) loss value in the corresponding database of minima. The graph shows how for the average test loss increases towards the bottom of the train loss landscape for  $N_{\text{data}} = 1000, 2000$ , as one would normally expect. For  $N_{\text{data}} = 100$  however, the average test loss seems to be decreasing again at the bottom of the train loss landscape.

#### 118 **S4. Example Disconnectivity Graphs for a Structural Glass-Former**

119 Two examples of disconnectivity graphs obtained for model structural glass-formers are shown in Figure S5  
120 for comparison with the loss function landscapes illustrated for neural networks. The system in question is a  
121 binary mixture of particles interacting via a Lennard-Jones potential (4) modelled with periodic boundary  
122 conditions and 60 or 256 particles in the supercell, described as BLJ<sub>60</sub> and BLJ<sub>256</sub>, respectively. Here,  
123 BLJ<sub>60</sub> contains 48 type A and 12 type B particles, while BLJ<sub>256</sub> contains 204 A and 52 B particles, and the  
124 results correspond to a number density of  $\sigma_{AA}^{-3}$ , where  $2^{1/6}\sigma_{AA}$  is the pair equilibrium separation for two  
125 A particles. The corresponding pair well depth is  $\epsilon_{AA}$ . Choosing  $\sigma_{AA} = 1$  and  $\epsilon_{AA} = 1$  defines a system  
126 of reduced units, and the additional parameters are  $\sigma_{AB} = 0.8$ ,  $\sigma_{BB} = 0.88$ ,  $\epsilon_{AB} = 1.5$ , and  $\epsilon_{BB} = 0.5$ .  
127 (15) The pairwise interactions were shifted and truncated according to the Stoddard-Ford scheme to assure  
128 continuous energy and first derivatives. The database for BLJ<sub>60</sub> contains over 11000 minima, and the  
129 database for BLJ<sub>256</sub> has 2500.

130 The potential energy landscapes in Figure S5 exhibit hierarchical structure, which appears to be common  
131 to other structural glasses (16–18), with numerous low-lying amorphous configurations separated by high  
132 barriers or order  $30k_B T_g$  for glass transition temperature  $T_g$ . Full details of the database construction and  
133 analysis can be found in the original reports (12–14).

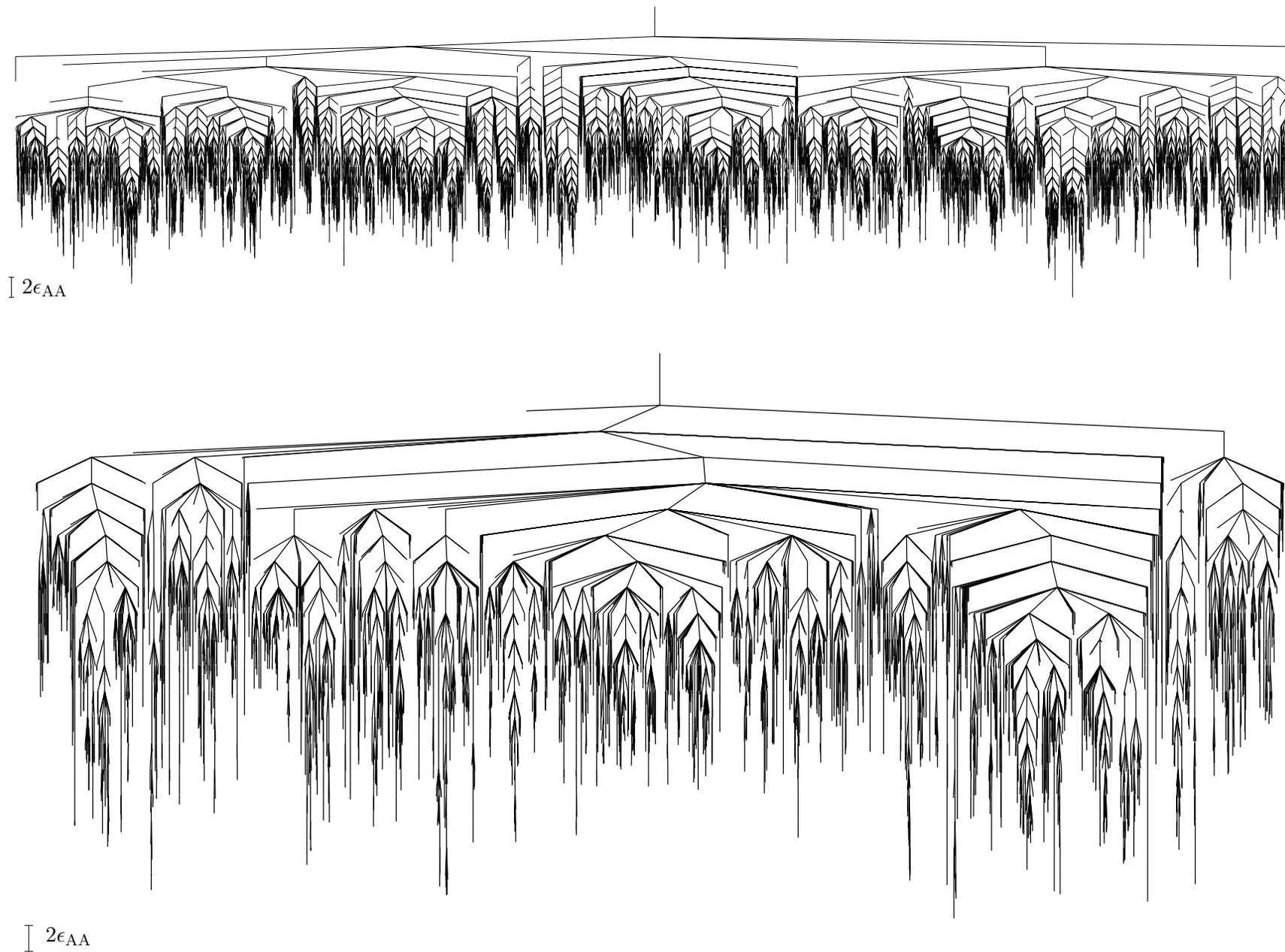


Fig. S5. Example disconnection graphs for binary Lennard-Jones systems containing 60 atoms, BLJ<sub>60</sub> (top), and 256 atoms, BLJ<sub>256</sub> (bottom), in periodically repeated supercells.(12–14)

134 **Movie S1. Graphical representation of the predictions for the D3 test set along the pathway**  
135 **shown in Fig. S3.**

## 136 **References**

- 137 1. A. J. Ballard, J. D. Stevenson, R. Das, and D. J. Wales. Energy landscapes for a machine learning  
138 application to series data. *J. Chem. Phys.*, 144(12):124119, 2016.
- 139 2. R. Das and D. J. Wales. Machine learning prediction for classification of outcomes in local minimisa-  
140 tion. *Chem. Phys. Lett.*, 667:158 – 164, 2017. .
- 141 3. D. J. Wales. Exploring energy landscapes. *Ann. Rev. Phys. Chem.*, 69:401–425, 2018. .
- 142 4. J. E. Jones and A. E. Ingham. On the calculation of certain crystal potential constants, and on the  
143 cubic crystal of least potential energy. *Proc. R. Soc. A*, 107:636–653, 1925.
- 144 5. B. M. Axilrod and E. Teller. Interaction of the van der waals type between three atoms. *J. Chem.*  
145 *Phys.*, 11:299, 1943.
- 146 6. P. G. Mezey. *Potential Energy Hypersurfaces*. Elsevier, Amsterdam, 1987.
- 147 7. D. J. Wales. *Energy Landscapes*. Cambridge University Press, Cambridge, 2003.
- 148 8. Stefano Martiniani, K. Julian Schrenk, Jacob D. Stevenson, David J. Wales, and Daan Frenkel. Turn-  
149 ing intractable counting into sampling: Computing the configurational entropy of three-dimensional  
150 jammed packings. *Phys. Rev. E*, 93:012906, 2016.
- 151 9. K. Swersky, J. Snoek, and R. P. Adams. Freeze-thaw bayesian optimization. *arXiv:1406.3896*  
152 *[stat.ML]*, 2014.
- 153 10. A. J. Ballard, R. Das, S. Martiniani, D. Mehta, L. Sagun, J. D. Stevenson, and D. J. Wales. Energy  
154 landscapes for machine learning. *Phys. Chem. Chem. Phys.*, 19:12585–12603, 2017. .
- 155 11. R. E. Leone and P. v. R. Schleyer. *Angew. Chem. Int. Ed. Engl.*, 9:860, 1970.
- 156 12. V. K. de Souza and D. J. Wales. Energy landscapes for diffusion: Analysis of cage-breaking processes  
157 (13 pages). *J. Chem. Phys.*, 129:164507, 2008.
- 158 13. V. K. de Souza and D. J. Wales. Connectivity in the potential energy landscape for binary lennard-  
159 jones systems (12 pages). *J. Chem. Phys.*, 130:194508, 2009.
- 160 14. V. K. de Souza, J. D. Stevenson, S. P. Niblett, J. D. Farrell, and D. J. Wales. Defining and quantifying  
161 frustration in the energy landscape: Applications to atomic and molecular clusters, biomolecules,  
162 jammed and glassy systems. *J. Chem. Phys.*, 146(12):124103, 2017. .
- 163 15. W. Kob and H. C. Andersen. Scaling behavior in the beta-relaxation regime of a supercooled lennard-  
164 jones mixture. *Phys. Rev. Lett.*, 73:1376–1379, 1994.
- 165 16. S. P. Niblett, V. K. de Souza, J. D. Stevenson, and D. J. Wales. Dynamics of a molecular glass former:  
166 Energy landscapes for diffusion in ortho-terphenyl. *J. Chem. Phys.*, 145(2):024505, 2016. .
- 167 17. S. P. Niblett, M. Biedermann, D. J. Wales, and V. K. de Souza. Pathways for diffusion in the potential  
168 energy landscape of the network glass former sio2. *The Journal of Chemical Physics*, 147(15):152726,  
169 2017. .
- 170 18. S. P. Niblett, V. K. de Souza, R. L. Jack, and D. J. Wales. Effects of random pinning on the potential  
171 energy landscape of a supercooled liquid. *The Journal of Chemical Physics*, 149(11):114503, 2018. .