## **Machine-learning for materials and physics discovery** through symbolic regression and kernel methods Stephen R. Xie, Shreyas Honrao, and Richard G. Hennig, University of Florida



**MPInterfaces** - High throughput framework for 2D materials

**VASPSol** - Ab initio methods for solid/liquid interfaces



Powered by **MPInterfaces & materialsweb** 



**GASP** - Genetic algorithm and machine learning for structure predictions

Open source available at <u>https://github.com/henniggroup</u>





### Data available at <a href="http://materialsweb.org">http://materialsweb.org</a>







# **Machine-learning for materials and physics discovery** through symbolic regression and kernel methods

Stephen R. Xie, Shreyas Honrao, and Richard G. Hennig, University of Florida

### **Machine Learning**

- Machine learning of energy landscapes using distribution functions
- Reduction of error by learning atomic energies with local RDF descriptors
- Learning of analytic equations to predict superconductivity using small data sets



Powered by **MPInterfaces & materialsweb** 



**GASP** - Genetic algorithm and machine learning for structure predictions

Open source available at <u>https://github.com/henniggroup</u>



**MPInterfaces** - High throughput framework for 2D materials



### Data available at <a href="http://materialsweb.org">http://materialsweb.org</a>







- GASP genetic algorithm: B. Revard, W. Tipton, A. Yesupenko, H. Lester,
- Machine learning of physics and materials: S. Honrao, S. Xie, B. Antonio
- Collaborators: Hao Li (UTAustin), Graeme Henkelman (UTAustin), Dallas R. Trinkle (UIUC)
- Machine learning of superconductivity: S. Xie, P. Hirschfeld, J. Hamlin, G. Stewart
- Financial support by NSF, DOE, NIST
- Computational resources: HiPerGator@UF, NSF XSEDE, and Google Cloud Platform





rhennig@ufl.edu http://hennig.mse.ufl.edu









# **Machine Learning in Materials Science**



Powered by MPInterfaces & materialsweb





# Part I: Exploration of Materials Energy Landscapes by Evolutionary/Genetic Algorithms

Powered by **MPInterfaces & materialsweb** 



rhennig@ufl.edu http://hennig.mse.ufl.edu





# **Genetic Algorithm Search for Crystalline Materials**





# Variable number of atoms and composition



# **Efficiency of Genetic Algorithm**

### **Efficiency compared to random search**

- 2.35 • Random search requires 2-3x more structure relaxations
- Genetic algorithm learns from previous structures

value Best 2.33

2.32

Powered by MPInterfaces & materialsweb



<u>https://github.com/henniggroup/gasp-python</u>

W. W. Tipton, RGH, J. Phys.: Cond. Matter 25, 495401 (2013) B. C. Revard, W. W. Tipton, A. Yesypenko, R. G. Hennig, PRB 93, 054117 (2016)







### **Problem:**

Naïve algorithm oversamples average compositions

### **Solutions**

- 1. Use larger endpoint structures
  - Works but expensive
- 2. Preferentially select parents with similar compositions
  - Needs metric for distance of structures in composition space





### https://github.com/henniggroup/gasp-python B. C. Revard and RGH, in preparation

# **Phase Diagram Searching**

### Metric for distance in composition space

- Express composition as a vector
- Use *L*<sub>1</sub> Norm to define distance:  $\bullet$

$$d_{XY} = \frac{1}{2} ||\mathbf{X} - \mathbf{Y}||_{1}$$
$$||\mathbf{A} - \mathbf{B}||_{1} = ||(1, -1, 0)||_{1} = |1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1| + |-1|$$

### **Fitness for structures**

- $f_{\rm comp} = 1 d$ Composition fitness
- $f_{\rm rel} = w_{\rm comp} f_{\rm comp} + (1 w_{\rm comp}) f_{\rm reg}$ **Relative fitness**

### **Fitness for structures**

- Sampling distribution improved but not uniform
- Use partial phase diagram searches if needed  $\bullet$

Powered by MPInterfaces & materialsweb







B. C. Revard and RGH, in preparation



Powered by **MPInterfaces & materialsweb** 





### https://github.com/henniggroup/gasp http://hennig.mse.ufl.edu

# Part II: Machine Learning of Energy Landscapes

*Powered by* MPInterfaces & materialsweb



<u>rhennig@ufl.edu</u> <u>http://hennig.mse.ufl.edu</u>



### **Machine Learning Regression**

- Takes a vector  $x \in \mathbb{R}^n$  as input and return a scalar y
- Must first construct a vector-based data representation of the crystal structure that encodes relevant physical information, *i.e.* chemical identity and position of the atoms

## Structure representation (features) should ideally fulfill three criteria

- (i) **Invariance** with respect to choice of unit cell and crystal symmetry
- (ii) Uniqueness, so no two different crystal structures have the same vector representation
- (iii)Continuity, such that the energy difference between two crystal structures with vector representations  $x_1$  and  $x_2$  goes to zero in the limit  $||x_1 - x_2|| \rightarrow 0$





# **Partial Radial Distribution Functions**

$$g_{AB}(r) = \frac{1}{N_A} \sum_{i=1}^{N_A} \sum_{j=1}^{\infty} \frac{1}{r^n} \exp\left[-\frac{\left(r - d_{ij}^{AB}\right)^2}{2\sigma_g^2}\right] \Theta(d_c)$$

- Captures primary distance dependence of bonds
- Criteria:
  - + Invariance
  - + Continuity
  - Uniqueness
- Cannot distinguish between *homometric structures*, *i.e.* structures of identical atoms that exhibit the same set of interatomic distances



S. Honrao, RGH at al., submitted (2018)







# **Partial Radial Distribution Functions**

$$g_{AB}(r) = \frac{1}{N_A} \sum_{i=1}^{N_A} \sum_{j=1}^{\infty} \frac{1}{r^n} \exp\left[-\frac{\left(r - d_{ij}^{AB}\right)^2}{2\sigma_g^2}\right] \Theta(d_c)$$

- Captures primary distance dependence of bonds
- Criteria:
  - + Invariance
  - + Continuity
  - Uniqueness
- Cannot distinguish between *homometric structures*, *i.e.* structures of identical atoms that exhibit the same set of interatomic distances



S. Honrao, RGH at al., submitted (2018)







### **Formation energies of Li-Ge structures**

- 14,168 Li-Ge structures from genetic algorithm search for novel Li-Ge compounds
- Includes relaxed and unrelaxed structures from DFT relaxations of structure search
- Formation energy relative to crystal structure of pure components

$$E_{\rm f} = E_{\rm tot} -$$

• E<sub>f</sub> is not simply counting bonds, sensitive to small changes in bonding character

### **Structure groups**

- Group the structures according to the basin of attraction  $\Rightarrow$  679 basin groups
- Splitting of data significantly reduces the correlation between testing and training set
- Provides more stringent and realistic test of the ML methods





 $X_{\mathrm{Li}}E_{\mathrm{Li}} - X_{\mathrm{Ge}}E_{\mathrm{Ge}}$ 

IPAM MLP Workshop I September 23-27, 2019 • UCLA

S. Honrao, RGH at al., submitted (2018)

### **Machine-learning models**

Use <u>kernel-ridge regression</u> (KRR), <u>ε-support vector regression</u> (SVR), and <u>neural networks</u>

### Input data preprocessing

- Feature scaling of components of input vector  $\mathbf{x}_i = g^i_{AB}(r)$  in training set to obtain zero mean and unit standard deviation
- Standardizing each set of components avoids the norm from being biased towards vector components with higher variance
- 30% of data for learning, 70% for testing
- 10-fold cross validation for learning











### Hyper parameter selection

- Determined from 10-fold cross validation
- $\epsilon$  for SVR: negligible changes for  $\epsilon < 10$  meV/atom, use  $\epsilon = 10$  meV/atom
- Cutoff distance varied from 5 40 Å, larger errors for 5 Å, select 10 Å

Algorithm	Kernel width	Average regularization parameter
KRR	21.2	0.015 (unit less)
SVR	54.3	10.5 meV/atom

• For neural network, 1 hidden layer and RELU function, 20 trials to estimate RMSE

Powered by **MPInterfaces & materialsweb** 



S. Honrao, RGH at al., Comp. Mater. Sci. 158, 414 (2019)



Algorithm	MAE	RMSE	R <sup>2</sup>
KRR	12.7	20.4	0.98
SVR	13.6	20.8	0.98
NN	12.8	20.2	0.98

- Similar prediction errors (meV/atom) for different machine-learning techniques
- NN most demanding
- SVR is computationally most efficient and provides best tradeoff between complexity and prediction error

## **Chemical accuracy for learning of energy landscape**



# Learning of Basin of Attractions from Unrelaxed Structures

Prediction of the relaxed energies (minima) from unrelaxed structures

Algorithm	MAE	RMS	<b>R</b> <sup>2</sup>
KRR	12.7	20.4	0.98
SVR	13.6	20.8	0.98
NN	12.8	20.2	0.98
KRR-min	11.8	20.3	0.98
SVR-min	13.4	20.9	0.98

- Similar accuracy for minima prediction
- Useful for screening of GA structures to avoid costly DFT relaxations

## ML learning of minima has similar accuracy as learning of energy landscape









- Often insufficient amount of data for ML
- Early stage of genetic algorithm searches:
  - Few dozen relaxed structure
  - Maybe 1,000 configurations
  - RMSE  $\approx$  35 meV/atom

### How can we improve the prediction error for small data sets?

## **Big Data in Materials Predictions?**





# **Data Augmentation - Use Local Descriptors and Information**

- Learn energy of individual atoms, separate ML models for each species Use local radial and angular distribution functions

$$g_i^{AB}(r) = \sum_{j=1}^{\infty} \frac{1}{r^2} \exp \left[ -\frac{\left(r - d_{ij}^{AB}\right)^2}{2\sigma_g^2} \right] f(d_{ij}^{AB})$$

Dataset: GASP run with Ni-Al EAM and Cd-Te Stillinger-Weber potential





$$q_i^{ABC}(x) = \sum_{j=1}^{\infty} \sum_{k=1}^{\infty} \exp\left[-\frac{\left(x - \cos\theta_{jik}\right)^{-}}{2\sigma_g^2}\right] f(d_{ij}^{AB}),$$

4,673 relaxed Cd-Te structures

rhennig@ufl.edu http://hennig.mse.ufl.edu



 $\sqrt{2}$ 





# **Prediction of Total Energy for Al-Ni**



### Significant improvement of prediction error for same amount of data.



• Angular terms do not change prediction error for Al-Ni, expected for EAM pair functional form







# **Prediction of Total Energy for Cd-Te**



Angular terms reduce prediction error for Cd-Te, expected for SW potential



Significant improvement of prediction error for same amount of data.





# Learning Curves



### Local RDF descriptor reduces prediction error even for small datasets.







# **Comparison to Other Descriptors for Cd-Te**

**Data Representation** Local RDF & ADF **Global RDF & ADF** Baseline model = mean  $E_{\rm f}$ **Coulomb matrix Orbital-field matrix Bag of bonds JARVIS-CFID AGNI fingerprints** 



RMSE (meV/atom)	MAE (meV/atom)
11	8
33	24
109	85
88	64
64	47
77	57
47	35
108	82

Importance of local RDF descriptors to capture change in chemical bonding



# Part III: Functional Form of the Superconducting **Critical Temperature from Machine Learning**

Stephen R. Xie, James Hamlin, Gregory R. Stewart, Peter J. Hirschfeld, Richard G. Hennig





# Motivation

- Vast space of possible superconducting materials
- Significant efforts to apply computational methods with theory to screen materials - 2015: Prediction and discovery of  $T_c$  in H<sub>3</sub>S at 200K and one megabar pressure

  - 2018: Prediction and discovery of T<sub>c</sub> in Lanthanum Hydrides between 215-260K (approaching room temperature)
- What are the correct material "descriptors" that reflect the underlying mechanism of superconductivity?
- Machine-learning of materials often suffers from a small-data problem. What methods can overcome this issue?
- Can we use machine learning to identify analytical relationships between descriptors and  $T_c$ ?



rhennig@ufl.edu http://hennig.mse.ufl.edu



- Training data from Table I by Allen & Dynes (1975)
  - 29 materials with parameters derived from tunneling
  - Parameters to use:  $\omega_{log}$ ,  $\lambda$ ,  $\mu^*$
  - Allen-Dynes: Most widely used Eq. to predict  $T_c$

PHYSICAL REVIEW B

VOLUME 12, NUMBER 3

### Transition temperature of strong-coupled superconductors reanalyzed

P. B. Allen\*

Department of Physics, State University of New York, Stony Brook, New York 11790

R. C. Dynes Bell Laboratories, Murray Hill, New Jersey 07974 (Received 20 January 1975)

$$T_{c} = \frac{f_{1}f_{2}\omega_{\log}}{1.20} \exp\left(-\frac{1.04(1+\lambda)}{\lambda - \mu^{*} - 0.62\,\lambda\mu^{*}}\right)$$

Powered by **MPInterfaces & materialsweb** 



rhennig@ufl.edu http://hennig.mse.ufl.edu

## **Training Data Set**



**1 AUGUST 1975** 

TABLE I. Parameters<sup>a</sup> of superconductors derived from tunneling measurements. The value of  $\mu^*$  is renormalized from previously reported values as described in the text.

Material	ω <sub>10g</sub> (K)	$\overline{\omega}_1$ (K)	$\overline{\omega}_2$ (K)	λ	$\omega_{ph}$ (K)	$\mu^*$ ( $\omega_{ph}$ )	$T_c$ (K)
Pb	56	60	65	1.55	110	0.105	7.2
In	68	79	89	0.805	179	0.097	3.40
Sn	99	110	121	0.72	209	0.092	3.75
Hg	29	38	49	1.6	162	0.098	4.19
TÌ	52	58	64	0.795	127	0.111	2.36
Та	132	140	148	0.69	228	0.093	4.48
a-Ga	55	77	101	1.62	291	0.095	8.56
β-Ga	87	108	129	0.97	285	0.092	5.90
$Tl_{0.9}Bi_{0.1}$	48	55	62	0.78	120	0.099	2.30
$Pb_{0.4}Tl_{0.6}$	48	56	62	1.15	121	0.094	4.60
$Pb_{0,6}Tl_{0,4}$	50	57	62	1.38	119	0.103	5.90
$Pb_{0.8}Tl_{0.2}$	50	56	61	1.53	116	0.101	6.80
$Pb_{0,6}Tl_{0,2}Bi_{0,2}$	48	53	58	1.81	112	0.111	7.26
$Pb_{0.9}Bi_{0.1}$	50	56	60	1.66	108	0.081	7.65
$Pb_{0.8}Bi_{0.2}$	46	52	57	1.88	109	0.093	7.95
$Pb_{0.7}Bi_{0.3}$	47	52	57	2.01	110	0.092	8.45
$Pb_{0.65}Bi_{0.35}$	45	50	55	2.13	110	0.093	8.95
$In_{0.9}Tl_{0.1}$	63	75	86	0.85	176	0.103	3.28
In <sub>0.73</sub> Tl <sub>0.27</sub>	55	67	77	0.93	166	0.110	3.36
In <sub>0.67</sub> Tl <sub>0.33</sub>	57	68	79	0.90	167	0.110	3.26
$In_{0.57}Tl_{0.43}$	53	64	<b>7</b> 4	0.85	165	0.117	2.60
$In_{0.5}Tl_{0.5}$	53	64	73	0.83	163	0.110	2.52
$In_{0.27}Tl_{0.73}$	42	53	63	1.09	151	0.094	3.64
$In_{0.17}Tl_{0.83}$	45	55	63	0.98	144	0.101	3.19
$In_{0.07}Tl_{0.93}$	49	56	63	0.89	131	0.107	2.77
$In_2Bi$	46	57	67	1.40	<b>174</b>	0.096	5.6
$Sb_2Tl_7$	37	48	58	1.43	<b>134</b>	0.102	5.2
$Bi_2Tl$	47	53	59	1.63	120	0.101	6.4
$a-\mathrm{Pb}_{0.45}\mathrm{Bi}_{0.55}$	29	38	47	2.59	128	0.116	7.0

<sup>a</sup>Tabulation of the data used to derive these parameters is available in J. M. Rowell, W. L. McMillan, and R. C. Dynes, J. Phys. Chem. Ref. Data (to be published).

> **IPAM MLP Workshop I** September 23-27, 2019 • UCLA

### $\eta (eV/Å^2)$ 2.4 1.3 2.2 1.4 1.24.9 2.12.01.1 1.6 2.0 2.12.32.22.3 2.42.41.4 1.4 1.4 1.3 1.3 1.4 1.31.3 1.6 1.6 2.12.1

# **Symbolic Regression Machine-Learning**

- Apply Sure Independence Screening and Sparsifying Operator (SISSO) method to generate predictive models from training data for 29 materials
  - Identifies *analytical relations* between a minimal set of descriptors and desired properties
  - Yields stable results with small training sets
- Compare equations by evaluating with testing data
  - 13 superconductors from literature including A15 phases
  - Measure interpolative and extrapolative capacity, i.e., transferability
  - Avoid overfitting, i.e., artificially low error in training data with high error in testing data





rhennig@ufl.edu http://hennig.mse.ufl.edu

- Primary features:  $\omega_{log}$ ,  $\lambda$ ,  $\mu^*$
- <u>Sure Independence Screening SIS</u>
  - 1. Expand feature space Φ by recursively applying and combining algebraic/functional operation
    - +, -, ×,  $\div$ , exp, log,  $\sqrt{1}$ ,  $\frac{1}{2}$ ,  $\frac{3}{3}$ , (sin, cos not used)
    - Respect units with dimensional reduction
  - 2. Rank features by their correlation magnitude (dot product of feature and  $T_c$ )
- <u>Sparsifying Operator SO</u>

  - Use  $\ell_0$ -norm regularized minimization to find sparse solution of linear equations - 1D solution is trivially the first-ranked feature
  - *n*-dimensional descriptors are used for classification models

Powered by **MPInterfaces & materialsweb** 



rhennig@ufl.edu http://hennig.mse.ufl.edu



• Feature: Quantity that is hypothesized to be relevant for describing target property,  $T_c$ 



Powered by MPInterfaces & materialsweb



rhennig@ufl.edu http://hennig.mse.ufl.edu









Powered by MPInterfaces & materialsweb



rhennig@ufl.edu http://hennig.mse.ufl.edu











Powered by MPInterfaces & materialsweb



rhennig@ufl.edu http://hennig.mse.ufl.edu











Powered by MPInterfaces & materialsweb



rhennig@ufl.edu http://hennig.mse.ufl.edu









Powered by **MPInterfaces & materialsweb** 



rhennig@ufl.edu http://hennig.mse.ufl.edu











Powered by MPInterfaces & materialsweb



rhennig@ufl.edu http://hennig.mse.ufl.edu





IPAM MLP Workshop I September 23-27, 2019 • UCLA

9 • • •





Powered by MPInterfaces & materialsweb



rhennig@ufl.edu http://hennig.mse.ufl.edu





IPAM MLP Workshop I September 23-27, 2019 • UCLA

. . . .





Powered by MPInterfaces & materialsweb



rhennig@ufl.edu http://hennig.mse.ufl.edu



λ				
exp, log,	,	-1, 2, 3,	9	$\sqrt{, \sqrt[3]{}}$
$exp(\mu^*)$	,	$\lambda^3$	9	$\sqrt{\mu^*}$
$\sqrt[3]{\mu^* + \lambda}$	, λ	$^{3\times}(\omega_{\log} \times \lambda)$	,	$exp(\lambda^3)$









rhennig@ufl.edu http://hennig.mse.ufl.edu



IPAM MLP Workshop I September 23-27, 2019 • UCLA

. . . .

9 • • •





### **Results in 3,414,094 analytic equations**

Powered by **MPInterfaces & materialsweb** 



rhennig@ufl.edu http://hennig.mse.ufl.edu



IPAM MLP Workshop I September 23-27, 2019 • UCLA

9 • • •



$$\begin{bmatrix} \Phi_0 \end{bmatrix} 3 \\ \begin{bmatrix} \Phi_1 \end{bmatrix} 34 \\ \begin{bmatrix} \Phi_2 \end{bmatrix} 1,342 \\ \begin{bmatrix} \Phi_3 \end{bmatrix} 3,414,094$$
  $\omega_{\log} \times \lambda$ ,  $\omega_{\log} \times \lambda$ ,  $\lambda^3 \times (\omega_{\log} \times \lambda)$ ,  $\lambda^$ 

### Select the equations with the highest linear correlation to $T_c$ , inner product > 0.5

Powered by MPInterfaces & materialsweb



<u>rhennig@ufl.edu</u> <u>http://hennig.mse.ufl.edu</u>





ening

# **Use Physical Constraints to Reduce Feature Space**

$$\begin{bmatrix} \Phi_0 \end{bmatrix} 3 \\ \begin{bmatrix} \Phi_1 \end{bmatrix} 34 \\ \begin{bmatrix} \Phi_2 \end{bmatrix} 1,342 \\ \begin{bmatrix} \Phi_3 \end{bmatrix} 3,414,094$$
  $\omega_{\log} \times \lambda$ ,  $\sqrt{\mu^*}$ ,  $\lambda^3 \times (\omega_{\log} \times \lambda)$ ,  $\lambda^3$   
 $\lambda^3 \times (\omega_{\log} \times \lambda)$ ,  $\lambda^3$   
 $\lambda^3 \times (\omega_{\log} \times \lambda) / (\lambda^3)$   
 $\lambda^3 \times (\omega_{\log} \times \lambda) / (\lambda^3)$   
Sure Independence Screet  
Dimensions  
 $15,886$   $\lambda \rightarrow 0$  Limit  
 $10,839$  Strictly Positive  
 $6,021$  Finite, Continuous, Real,

Powered by **MPInterfaces & materialsweb** 



rhennig@ufl.edu http://hennig.mse.ufl.edu





### ening

### Monotonic

### Tools available at <u>https://github.com/henniggroup/symbolic-regression-utilities</u>

- Fit multiplicative factor to features, e.g.,  $y' = C \frac{\lambda^4 \omega_{\log}}{\lambda^2 + \sqrt{\mu^*}}$
- Evaluate root-mean square (relative) error with leave-one-out cross-validation
- Figure shows distribution of resulting functions for  $\mu^* = 0.1$
- Narrow distribution of training and testing data
- Data for large and small  $\lambda$  would be helpful

Powered by MPInterfaces & materialsweb







# **Best Model Performance**



![](_page_43_Picture_2.jpeg)

# **Best Model Performance**

![](_page_44_Figure_1.jpeg)

Outliers (MgB<sub>2</sub>, NbS<sub>2</sub>) indicate importance of anisotropic electron-phonon coupling.

![](_page_44_Picture_3.jpeg)

![](_page_44_Picture_4.jpeg)

# **Dimensionality and Complexity**

Model	CV-RMSE (K) Training	RMSE (K) Testing	Materials Parameters	Numerical Coefficients
SISSO	0.25	3.2	3	1
mod. McMillan	0.92		3	4
Allen-Dynes	0.30		4	7

### Machine-learned model has small RMSE and low computational complexity.

*Powered by* MPInterfaces & materialsweb

![](_page_45_Picture_4.jpeg)

<u>rhennig@ufl.edu</u> <u>http://hennig.mse.ufl.edu</u>

![](_page_45_Picture_6.jpeg)

### **Can the machine-learning identify relevant materials parameters?**

- Fit models with up to 7 materials parameters for 29 materials
- Best model:

$$T_c^{\rm SISSO} = 0.09525 \frac{\lambda^4 \omega_{\log}}{\lambda^3 + \sqrt{\mu^*}}$$

• Second best model:

$$T_c = -0.059 \left(\omega_2 - \omega_1 - \frac{\omega_2}{\lambda}\right) \frac{\lambda^3}{\sqrt[3]{\lambda}}$$

• Adding parameters beyond  $\omega_{\log}$ ,  $\lambda$ ,  $\mu^*$  does not improve description

### Machine-learning identifies most relevant materials parameters in agreement with McMillan and Allen & Dynes

Powered by MPInterfaces & materialsweb

![](_page_46_Picture_10.jpeg)

rhennig@ufl.edu http://hennig.mse.ufl.edu

V

![](_page_46_Picture_12.jpeg)

### CV-RMSE = 0.25 K

### CV-RMSE = 0.27K

Parameters<sup>a</sup> of superconductors derived from tunneling measurements.  $\mu^*$  is renormalized from previously reported values as described in the text

•		-	-					
Material	ω <sub>10g</sub> (K)	<u>ω</u> <sub>1</sub> (K)	<u>ω</u> 2 (K)	λ	ω <sub>ph</sub> (K)	$\mu^*(\omega_{ph})$	<i>T</i> <sub>c</sub> (K)	$\eta$ (eV/Å <sup>2</sup> )
Pb	56	60	65	1.55	110	0.105	7.2	2.4
In	68	79	89	0.805	179	0.097	3.40	1.3
Sn	99	110	121	0.72	209	0.092	3.75	2.2
Hg	29	38	49	1.6	162	0.098	4.19	1.4
TÌ	52	58	64	0.795	127	0.111	2.36	1.2
Та	132	140	148	0.69	228	0.093	4.48	4.9
a-Ga	55	77	101	1.62	291	0.095	8.56	2.1
β-Ga	87	108	129	0.97	285	0.092	5.90	2.0
$Tl_{0.9}Bi_{0.1}$	48	55	62	0.78	120	0.099	2.30	1.1
$Pb_{0.4}Tl_{0.6}$	48	56	62	1.15	121	0.094	4.60	1.6
$Pb_{0,6}Tl_{0,4}$	50	57	62	1.38	119	0.103	5.90	2.0
$Pb_{0.8}Tl_{0.2}$	50	56	61	1.53	116	0.101	6.80	2.1
$Pb_{0} {}_{6}Tl_{0} {}_{2}Bi_{0} {}_{2}$	48	53	58	1.81	112	0.111	7.26	2.3
$Pb_{0.9}Bi_{0.1}$	50	56	60	1.66	108	0.081	7.65	2.2
$Pb_{0.8}Bi_{0.2}$	46	52	57	1.88	109	0.093	7.95	2.3
$Pb_{0,7}Bi_{0,3}$	47	52	57	2.01	110	0.092	8.45	2.4
$Pb_{0,65}Bi_{0,35}$	45	50	55	2.13	110	0.093	8.95	2.4
$In_{0.9}Tl_{0.1}$	63	75	86	0.85	176	0.103	3.28	1.4
$In_{0.73}Tl_{0.27}$	55	67	77	0.93	166	0.110	3.36	1.4
In <sub>0.67</sub> Tl <sub>0.33</sub>	57	68	79	0.90	167	0.110	3.26	1.4
$In_{0.57}Tl_{0.43}$	53	<b>64</b>	<b>7</b> 4	0.85	165	0.117	2.60	1.3
$In_{0.5}Tl_{0.5}$	53	64	73	0.83	163	0.110	2.52	1.3
$In_{0.27}Tl_{0.73}$	42	53	63	1.09	151	0.094	3.64	1.4
$In_{0,17}Tl_{0,83}$	45	55	63	0.98	144	0.101	3.19	1.3
$In_{0.07}Tl_{0.93}$	49	56	63	0.89	131	0.107	2.77	1.3
In <sub>2</sub> Bi	46	57	67	1.40	174	0.096	5.6	1.6
$Sb_2T1_7$	37	48	58	1.43	134	0.102	5.2	1.6
Bi <sub>2</sub> Tl	47	53	59	1.63	120	0.101	6.4	2.1
$a - \mathbf{Pb}_{0.45} \mathbf{Bi}_{0.55}$	29	38	47	2.59	128	0.116	7.0	2.1

ion of the data used to derive these parameters is available in J. M. Rowell, W. L. McMilla and R. C. Dynes, J. Phys. Chem. Ref. Data (to be published)

- Can increasing the number of numerical coefficients improve description Include multiplicative and additive numerical coefficient to every materials parameter
- Best model:

$$T_{c} = \omega_{\log} \left( \frac{0.233 - 0.0170\lambda}{1.28\mu^{*} + 0.00784} + (0.791\lambda - 1.408)^{3} \right) \left( 0.0655\lambda + 0.00530 - \frac{0.000780}{1.206\mu^{*} - 0.0725} - \frac{0.000780}{1.206\mu^{*} - 0.0725} \right)$$

$$CV-RMSE = 0.19 \text{ K}$$

 Significant improvement from 0.25 to 0.19 K, however, at cost of significant more complexity (increase from 1 to 11 numerical coefficients) and reduced physical interpretability (high power of  $\lambda$ )

![](_page_47_Picture_7.jpeg)

![](_page_47_Picture_9.jpeg)

rhennig@ufl.edu http://hennig.mse.ufl.edu

![](_page_47_Picture_12.jpeg)

![](_page_47_Picture_13.jpeg)

# Conclusions

- Machine-learning of analytic equation with fewer parameters
- Identification of relevant physical parameters
- Use of analytic expressions and physical constraints can help overcome small-data problem
- Predict known superconductors of
- same type as the original Allen-Dynes dataset Anomalous outliers suggests need for new descriptors anisotropy of the electron-phonon interaction

$$T_c^{\rm SISSO} = 0.09525 \frac{\lambda^4 \omega_{\log}}{\lambda^3 + \sqrt{\mu^*}}$$

Powered by **MPInterfaces & materialsweb** 

![](_page_48_Picture_8.jpeg)

![](_page_48_Picture_10.jpeg)

- (b) Testing and Extrapolation 40 **Training Data** Testing data  $\omega_{log}$ ,  $\lambda$  from tunneling •  $\omega_{\text{log}}$  from  $\omega_{\text{D}}$ ,  $\lambda$  from  $\rho(T)$ 30  $\star \omega_{\rm los}$  from  $\omega_{\rm D}$ ,  $\lambda$  from calc  $T_c$  predicted (K) Nb<sub>2</sub>Sn 20-Nb<sub>2</sub>Ge CaC Nb<sub>2</sub>Z LuNi<sub>2</sub>B<sub>2</sub>C Nb 10 La<sub>3</sub>Ni<sub>2</sub>B<sub>2</sub>N<sub>3</sub> 30 20 10 ()  $T_{c}$  from experiment (K) IPAM MLP Workshop I
- rhennig@ufl.edu http://hennig.mse.ufl.edu

September 23-27, 2019 • UCLA

# **Machine-learning for materials and physics discovery** through symbolic regression and kernel methods

Stephen R. Xie, Shreyas Honrao, and Richard G. Hennig, University of Florida

### **Search for materials**

### **MPInterfaces** - High throughput framework for 2D materials

- Structure prediction by genetic algorithms
- Machine learning of energy landscapes using distribution functions
- Learning of Allen-Dynes equation for  $T_c$
- Use of analytic equations and physical constraints to overcome small data problem

![](_page_49_Figure_8.jpeg)

**MPInterfaces & materialsweb** 

Powered by

and machine learning for structure predictions

**GASP** - Genetic algorithm

Open source available at <u>https://github.com/henniggroup</u>

![](_page_49_Picture_11.jpeg)

![](_page_49_Figure_12.jpeg)

### Data available at <a href="http://materialsweb.org">http://materialsweb.org</a>

![](_page_49_Picture_14.jpeg)

![](_page_49_Picture_16.jpeg)

![](_page_49_Picture_17.jpeg)