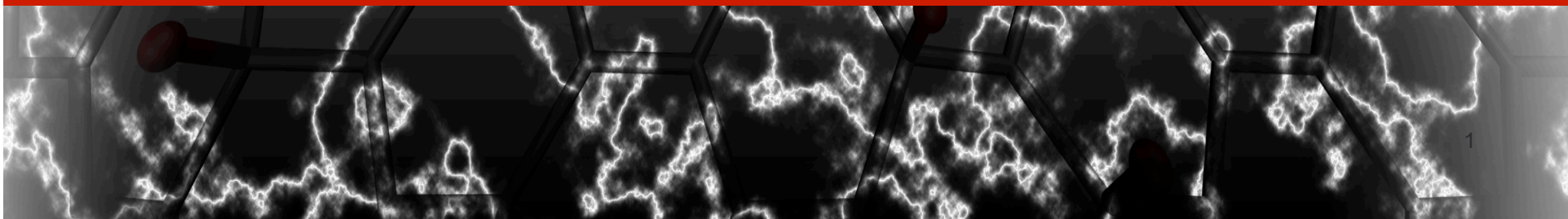


# How data curation has potential to help you do science

**Austin Sendek, E. Dogus Cubuk, Gowoon Cheon, Qian Yang,  
Yi Cui, Evan Reed\***

*\*Department of Materials Science and Engineering  
Stanford University*



# THE PLAN

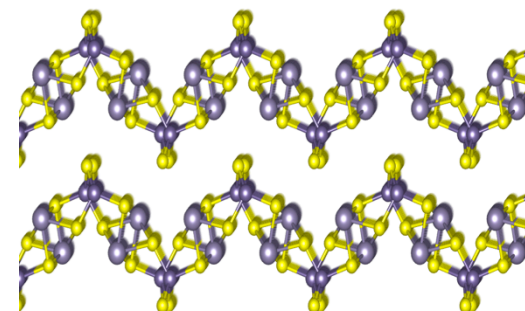


- Practical data curation: Concrete examples in 2D
- Why collect and curate data?
  - A machine learning example with experimentally measured data
  - How to know when to believe machine learning
  - How much data does one need?

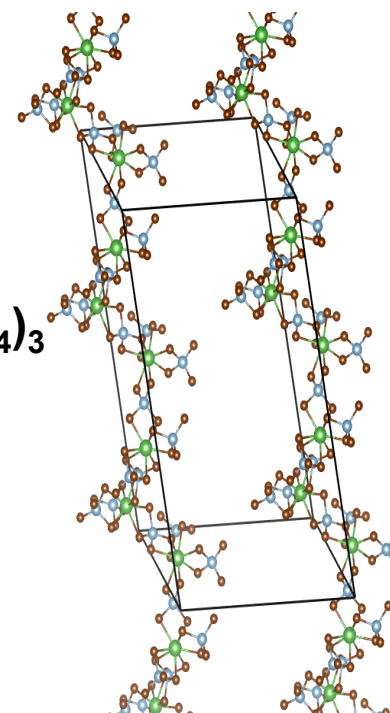
# We discover new 2D and 1D materials

- Data mining of public databases leads to the discovery of:
  - **1173** 2D layered materials
  - **325** materials with piezoelectric monolayers
  - **98** bulk vertical lattice-commensurate heterostructures
  - **487** 1D molecular wires

**SnGeS<sub>3</sub>**



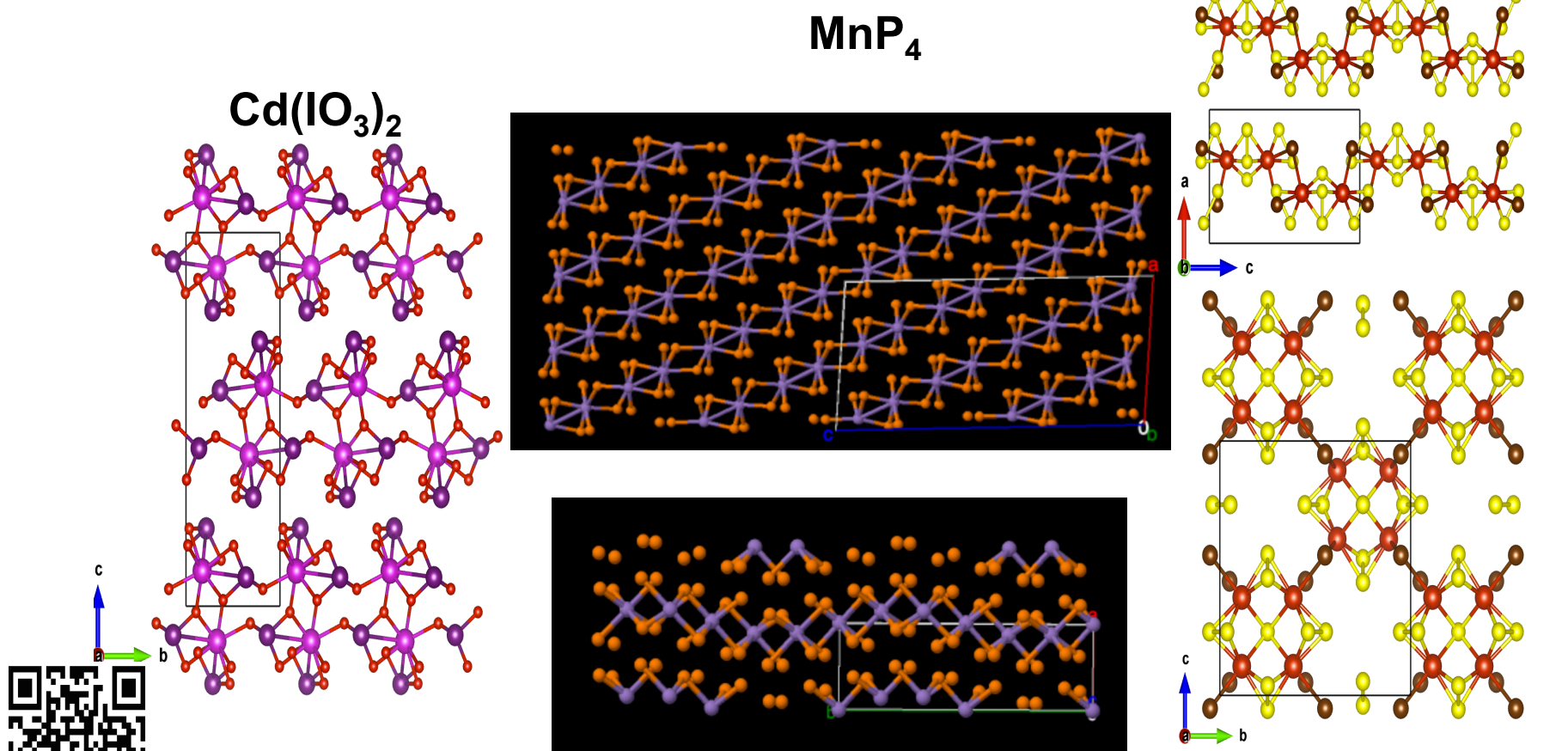
**La(AlBr<sub>4</sub>)<sub>3</sub>**



G.Cheon et al., Data mining for new two- and one-dimensional weakly bonded solids and lattice-commensurate heterostructures, *Nano Letters* (2017)

# We compile a genome of 1173 2D materials

- Diverse spectrum of layered materials
- Materials Project IDs of all layered materials available in Supporting Information



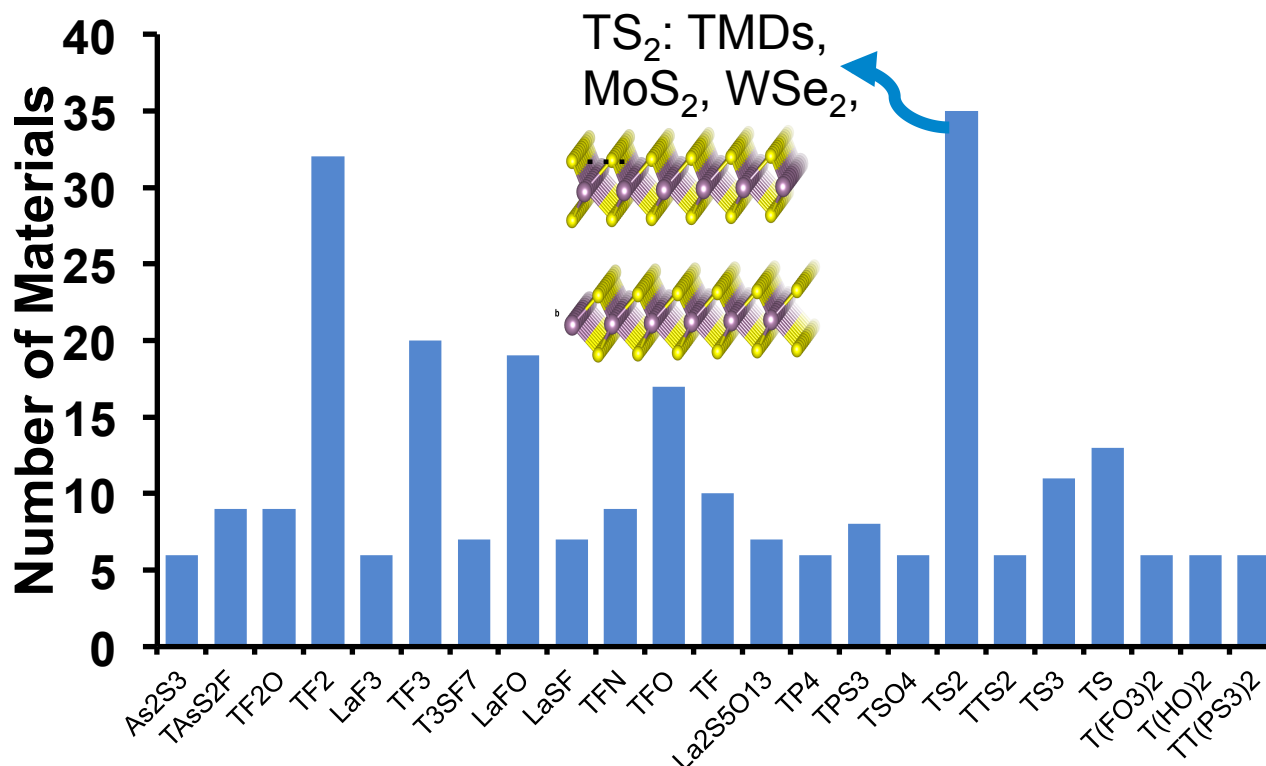
G.Cheon et al., *Nano Letters* (2017)



# We compile a genome of 1173 2D materials

- 1173 weakly bonded layered materials identified, lots of new candidates!
- 23 families of similar chemical compositions (>5 materials), but >80% don't belong to a family

## Families of 2D Materials

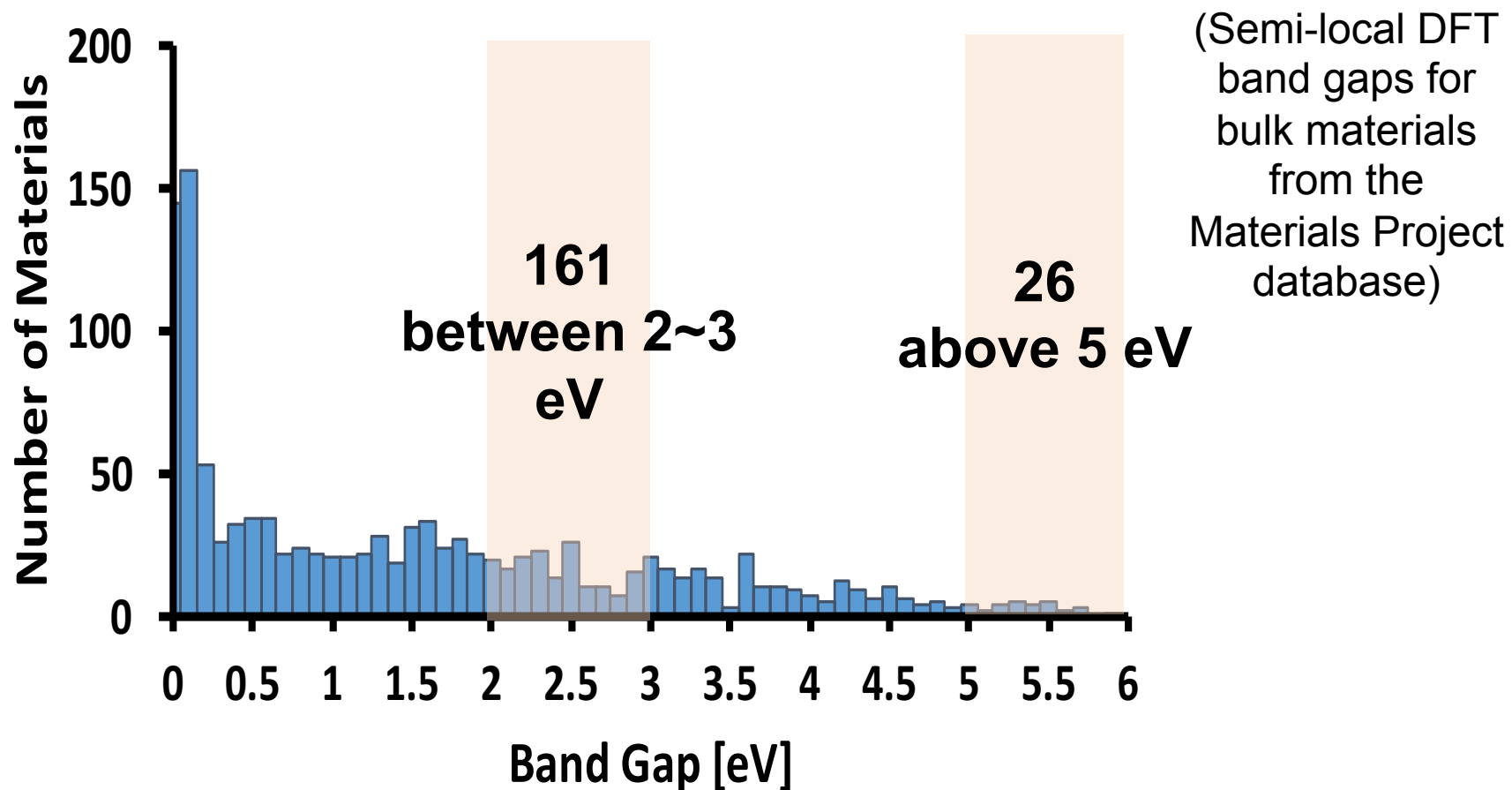


Ac: actinides  
As: large pnictogens (As, Sb, Bi)  
F: halogens  
S: chalcogens excluding O  
La: lanthanides  
T: transition metals



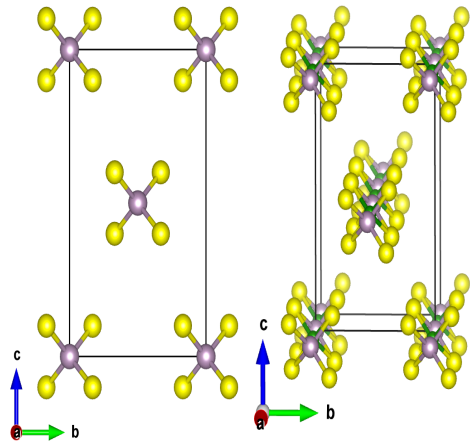
# We find a wide spectrum of 2D material band gaps

## Band Gap Distribution (2D Materials)

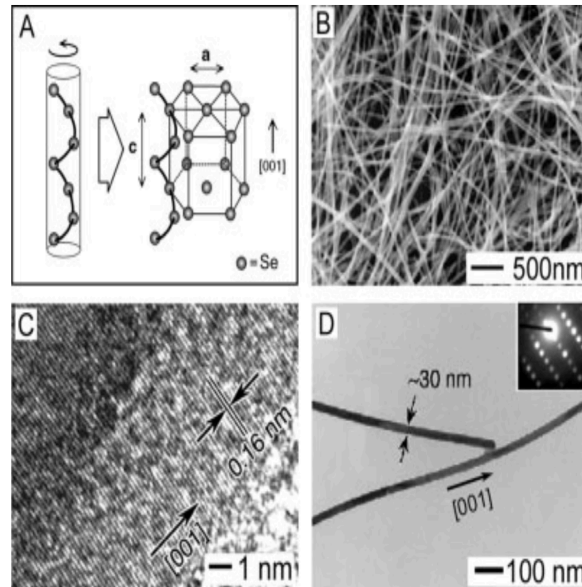


# We discover 487 1D materials

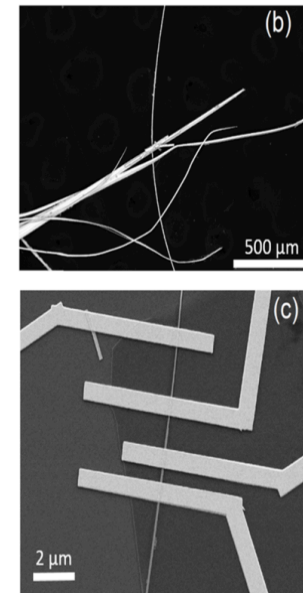
- Our algorithm can find 1D dimensional subunits



BPS<sub>4</sub>, a material found in this work



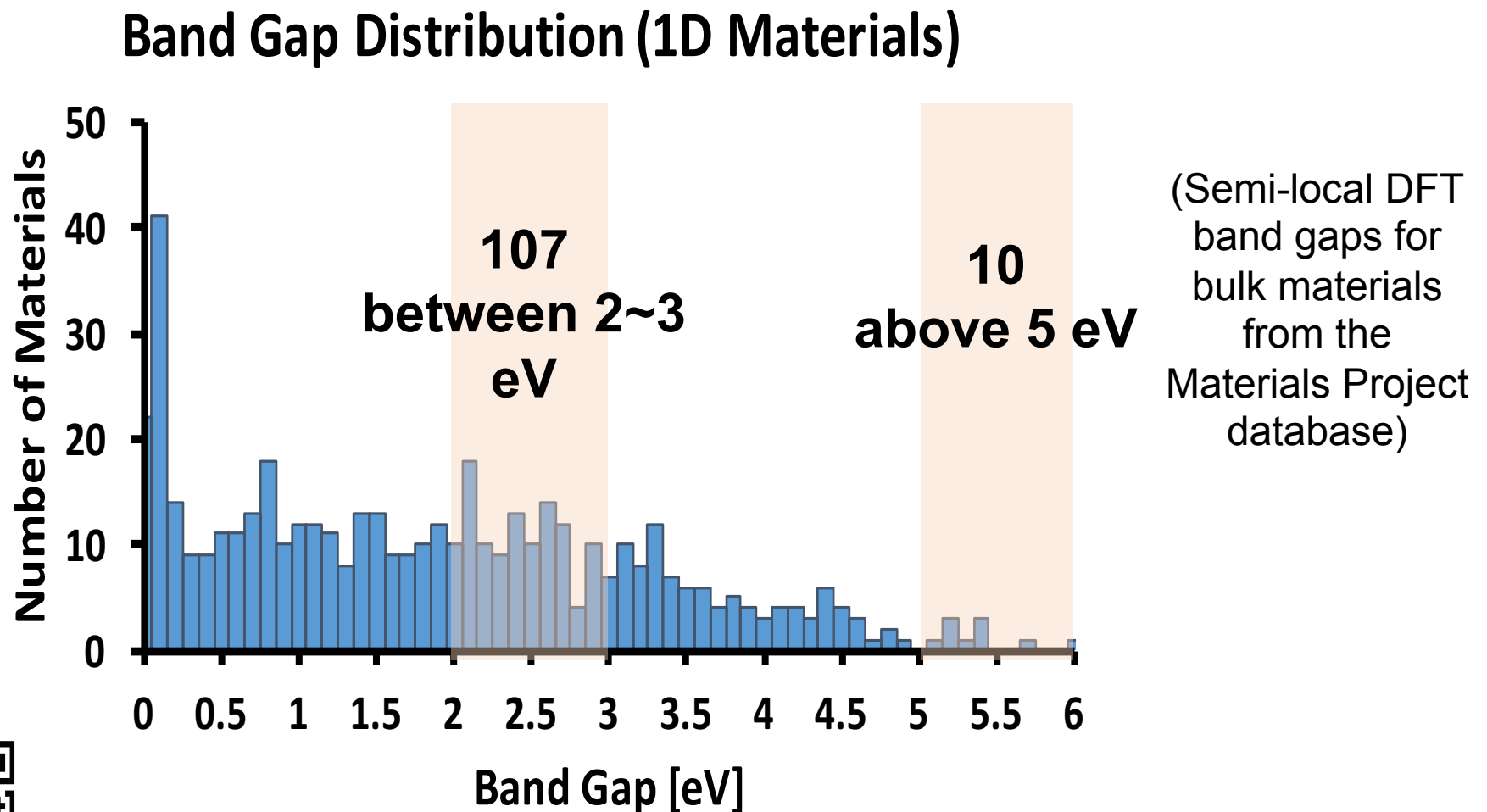
Chain-like structure of t-Se that grow into nanowires.  
Xia et al., *Adv. Materials* (2003)



Quasi-1D TaSe<sub>3</sub> low-noise nanowire devices  
Liu et al., *Nano Lett.* (2017)

- Some inorganic ‘molecular wires’ have been predicted to possess structural stability and versatile material properties, **but only ~20 known**

# We find a wide spectrum of 1D material band gaps



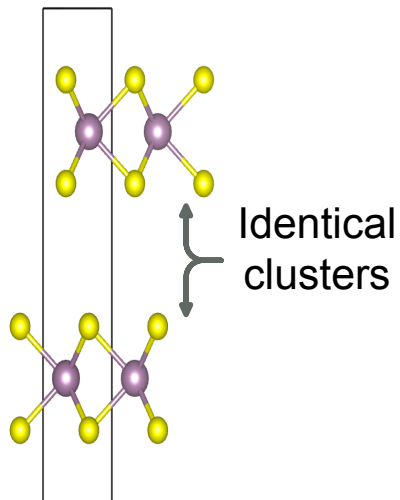


# We discover lattice-commensurate vertical heterostructures

- We discover **intrinsic, lattice-commensurate heterostructures** that preclude the need for artificial stacking:

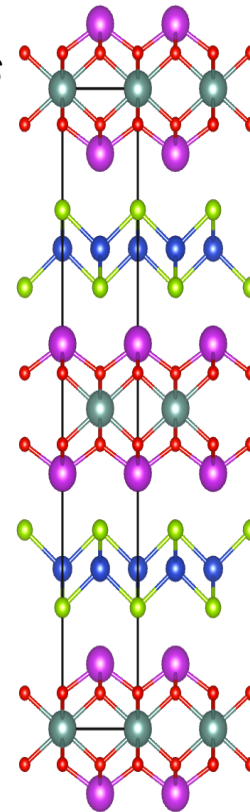
## Not a heterostructure

MoS<sub>2</sub>

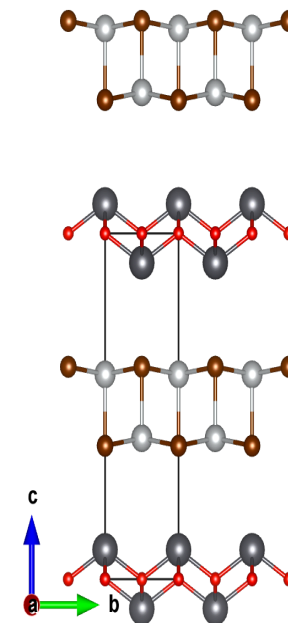


**Heterostructures:** YCu<sub>2</sub>Bi<sub>2</sub>(SeO<sub>2</sub>)<sub>2</sub>  
different  
*chemical compositions*  
or  
*number of atoms*  
in each layer

**Different clusters**



AgPbBrO

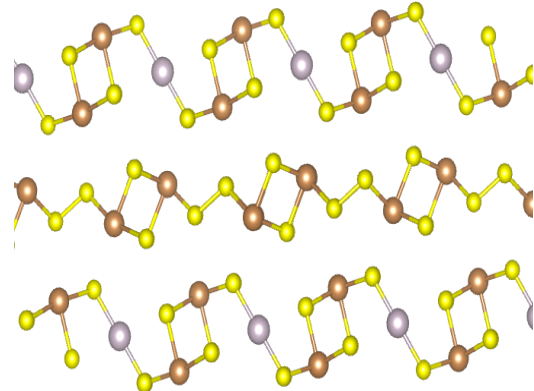


G.Cheon et al., *Nano Letters* (2017)

# We discover lattice-commensurate heterostructures

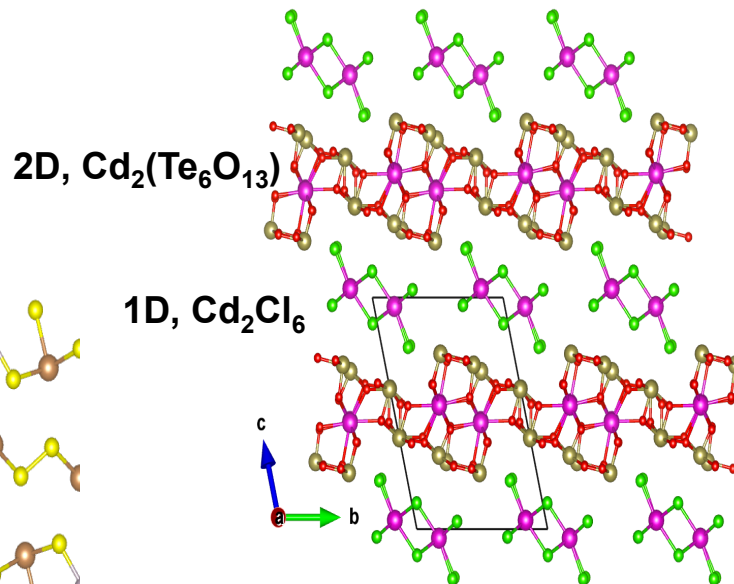
- We identify 98 lattice-commensurate heterostructures:
- Experimentally reported in bulk crystals

**HgSb<sub>4</sub>S<sub>8</sub>(Livingstonite)**



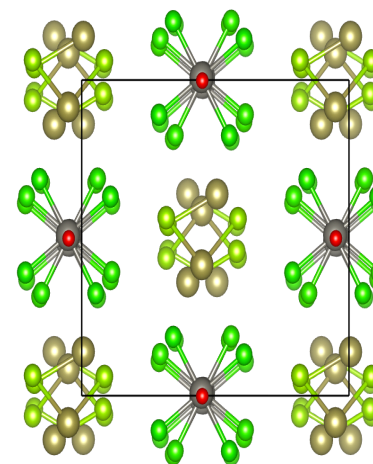
**2D**

**Cd<sub>4</sub>Te<sub>6</sub>Cl<sub>6</sub>O<sub>13</sub>**



**2D+1D**

**Te<sub>3</sub>W<sub>2</sub>Se<sub>4</sub>(Cl<sub>4</sub>O)<sub>2</sub>**



**1D**



# SOME PRACTICAL CURATION APPROACHES



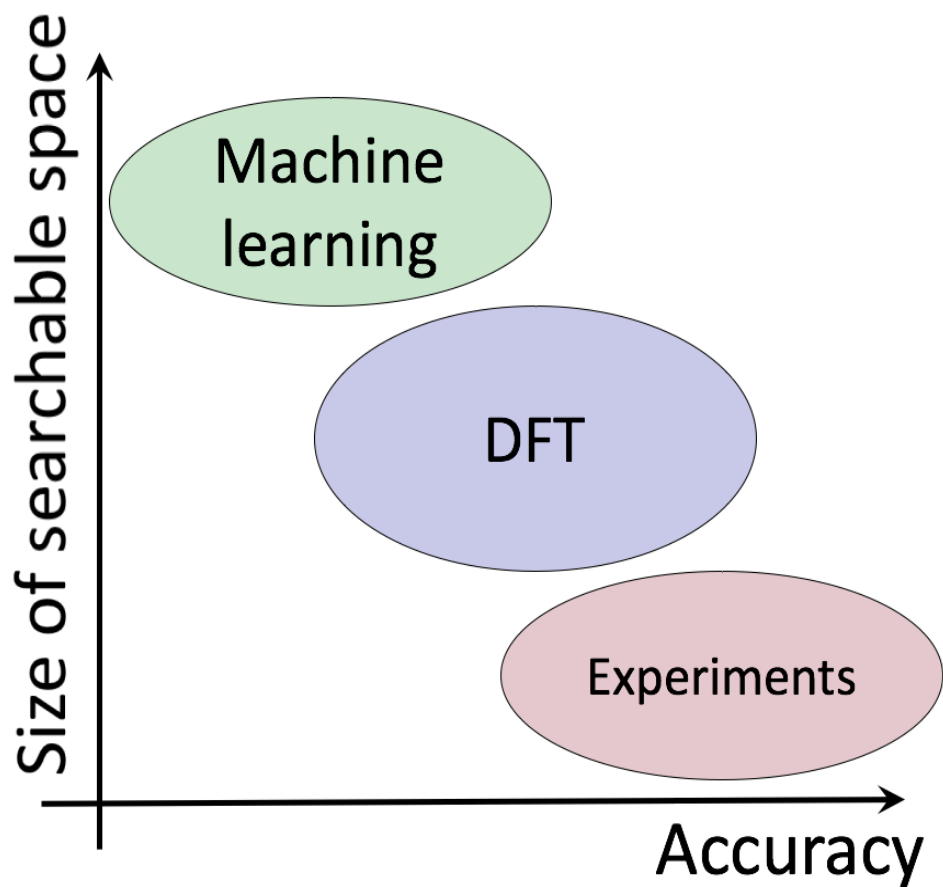
## Lower effort:

- Supporting information of publications
- File on your webpage
- File on Materials Data Facility (NIST)

## More effort, but broader utility:

- Work with others to fold into existing databases:
  - NSF 2DCC at Penn State (Vin Crespi, Richard Hennig, et al)
  - Jarvis at NIST (Francesca Tavazza, et al)
  - Materialsweb.org (Richard Hennig)
  - DOE's Materials Project (Kristin Persson et al)

MACHINE LEARNING HAS THE POTENTIAL TO FILL A GAP,  
ENABLING SEARCHES OF LARGE SPACES OF MATERIALS





# WE COLLECT 39 EXPERIMENTAL MEASUREMENTS OF LI ION CONDUCTIVITY FOR SOLIDS

Composition	RT bulk ionic conductivity (S cm <sup>-1</sup> )
LiLa(TiO <sub>3</sub> ) <sub>2</sub>	1 × 10 <sup>-3</sup>
Li <sub>9.81</sub> Sn <sub>0.81</sub> P <sub>2.19</sub> S <sub>12</sub>	5.5 × 10 <sup>-3</sup>
Li <sub>10</sub> Ge(PS <sub>6</sub> ) <sub>2</sub>	1.4 × 10 <sup>-2</sup>
Li <sub>10.35</sub> Si <sub>1.35</sub> P <sub>1.65</sub> S <sub>12</sub>	6.5 × 10 <sup>-3</sup>
Li <sub>14</sub> ZnGe <sub>4</sub> O <sub>16</sub> (2)	1.0 × 10 <sup>-6</sup>
Li <sub>2</sub> Ca(NH) <sub>2</sub>	6.4 × 10 <sup>-6</sup>
Li <sub>2</sub> Ge <sub>7</sub> O <sub>15</sub>	5.0 × 10 <sup>-6</sup>
Li <sub>2</sub> NH	2.5 × 10 <sup>-4</sup>
Li <sub>2</sub> S	1.0 × 10 <sup>-13</sup>
Li <sub>13.6</sub> Si <sub>2.8</sub> S <sub>1.2</sub> O <sub>16</sub>	6.0 × 10 <sup>-7</sup>
Li <sub>14</sub> Ge <sub>2</sub> V <sub>2</sub> O <sub>16</sub>	7.0 × 10 <sup>-5</sup>
Li <sub>15</sub> Ge <sub>3</sub> V <sub>2</sub> O <sub>4</sub>	6.03 × 10 <sup>-6</sup>
Li <sub>14.8</sub> Ge <sub>3.4</sub> W <sub>0.6</sub> O <sub>4</sub>	4.0 × 10 <sup>-5</sup>
Li <sub>3</sub> Fe <sub>2</sub> P <sub>3</sub> O <sub>12</sub>	1.0 × 10 <sup>-7</sup>
Li <sub>3</sub> N	5.75 × 10 <sup>-4</sup>
Li <sub>3</sub> P	1.0 × 10 <sup>-3</sup>
γ-Li <sub>3</sub> PS <sub>4</sub>	3.0 × 10 <sup>-7</sup>
Li <sub>3</sub> Sc <sub>2</sub> P <sub>3</sub> O <sub>12</sub>	1.0 × 10 <sup>-10</sup>
β <sub>11</sub> -Li <sub>3</sub> VO <sub>4</sub>	4.4 × 10 <sup>-8</sup>
Li <sub>4</sub> B <sub>7</sub> O <sub>12</sub> Cl	1.0 × 10 <sup>-7</sup>
Li <sub>4</sub> BN <sub>3</sub> H <sub>10</sub>	2.0 × 10 <sup>-4</sup>
γ-Li <sub>4</sub> GeO <sub>4</sub>	3.1 × 10 <sup>-12</sup>
Li <sub>4</sub> SiO <sub>4</sub>	2.4 × 10 <sup>-10</sup>
Li <sub>5</sub> La <sub>3</sub> Bi <sub>2</sub> O <sub>12</sub>	2.0 × 10 <sup>-5</sup>
Li <sub>5</sub> La <sub>3</sub> Nb <sub>2</sub> O <sub>12</sub>	8.0 × 10 <sup>-6</sup>
Li <sub>5</sub> La <sub>3</sub> Ta <sub>2</sub> O <sub>12</sub>	1.5 × 10 <sup>-6</sup>
Li <sub>5</sub> Ni <sub>2</sub>	1.5 × 10 <sup>-7</sup>
Li <sub>6</sub> BaLa <sub>2</sub> Ta <sub>2</sub> O <sub>12</sub>	4.0 × 10 <sup>-5</sup>
Li <sub>6</sub> FeCl <sub>8</sub>	1.0 × 10 <sup>-4</sup>
Li <sub>6</sub> NBr <sub>3</sub>	1.5 × 10 <sup>-7</sup>
Li <sub>6</sub> SrLa <sub>2</sub> Ta <sub>2</sub> O <sub>12</sub>	7.0 × 10 <sup>-6</sup>
Li <sub>7</sub> La <sub>3</sub> Zr <sub>2</sub> O <sub>12</sub>	3.5 × 10 <sup>-4</sup>
Li <sub>7</sub> P <sub>3</sub> S <sub>11</sub>	4.1 × 10 <sup>-3</sup>
LiAlH <sub>4</sub>	2.0 × 10 <sup>-9</sup>
LiAlSiO <sub>4</sub>	1.4 × 10 <sup>-5</sup>
LiBH <sub>4</sub>	2.0 × 10 <sup>-8</sup>
LiI	1.0 × 10 <sup>-6</sup>
LiNH <sub>2</sub>	4.0 × 10 <sup>-10</sup>
α'-LiZr <sub>2</sub> P <sub>3</sub> O <sub>12</sub>	5.0 × 10 <sup>-8</sup>

- We adopt a binary classification strategy with a 10<sup>-4</sup> S/cm boundary, motivated by engineering requirements
- Training set includes 8 “good” conductors, 31 “bad” conductors

A. D. Sendek, E. J. Reed, et al, Energy and Environmental Science (2017).

# WE DRAW ON WISDOM/PROPOSALS IN THE LITERATURE FOR READILY COMPUTABLE FEATURES (NO DFT!)

	Feature	Pearson correlation coefficient
1	Volume per atom <sup>a</sup>	0.20
2	Standard deviation in Li neighbour count	0.22
3	Standard deviation in Li bond ionicity	-0.04
4	Li bond ionicity <sup>a</sup>	-0.18
5	Li neighbour count <sup>a</sup>	-0.19
6	<i>Li-Li bonds per Li<sup>a</sup></i>	0.06
7	<i>Bond ionicity of sublattice<sup>a</sup></i>	-0.28
8	Sublattice neighbour count <sup>a</sup>	-0.13
9	<i>Anion framework coordination<sup>a</sup></i>	-0.06
10	Minimum anion-anion separation distance <sup>a</sup> (Å)	0.09
11	Volume per anion (Å <sup>3</sup> )	-0.01
12	<i>Minimum Li-anion separation distance<sup>a</sup> (Å)</i>	0.20
13	<i>Minimum Li-Li separation distance<sup>a</sup> (Å)</i>	-0.10
14	Electronegativity of sublattice <sup>a</sup>	-0.16
15	Packing fraction of full crystal	0.16
16	Packing fraction of sublattice	0.19
17	Straight-line path width <sup>a</sup> (Å)	0.07
18	Straight-line path electronegativity <sup>a</sup>	-0.29
19	Ratio of features (4) and (7)	-0.03
20	Ratio of features (5) and (8)	-0.18
	Constant term	—

No single feature has strong correlation with ionic conductivity across the broad spectrum of 39 materials

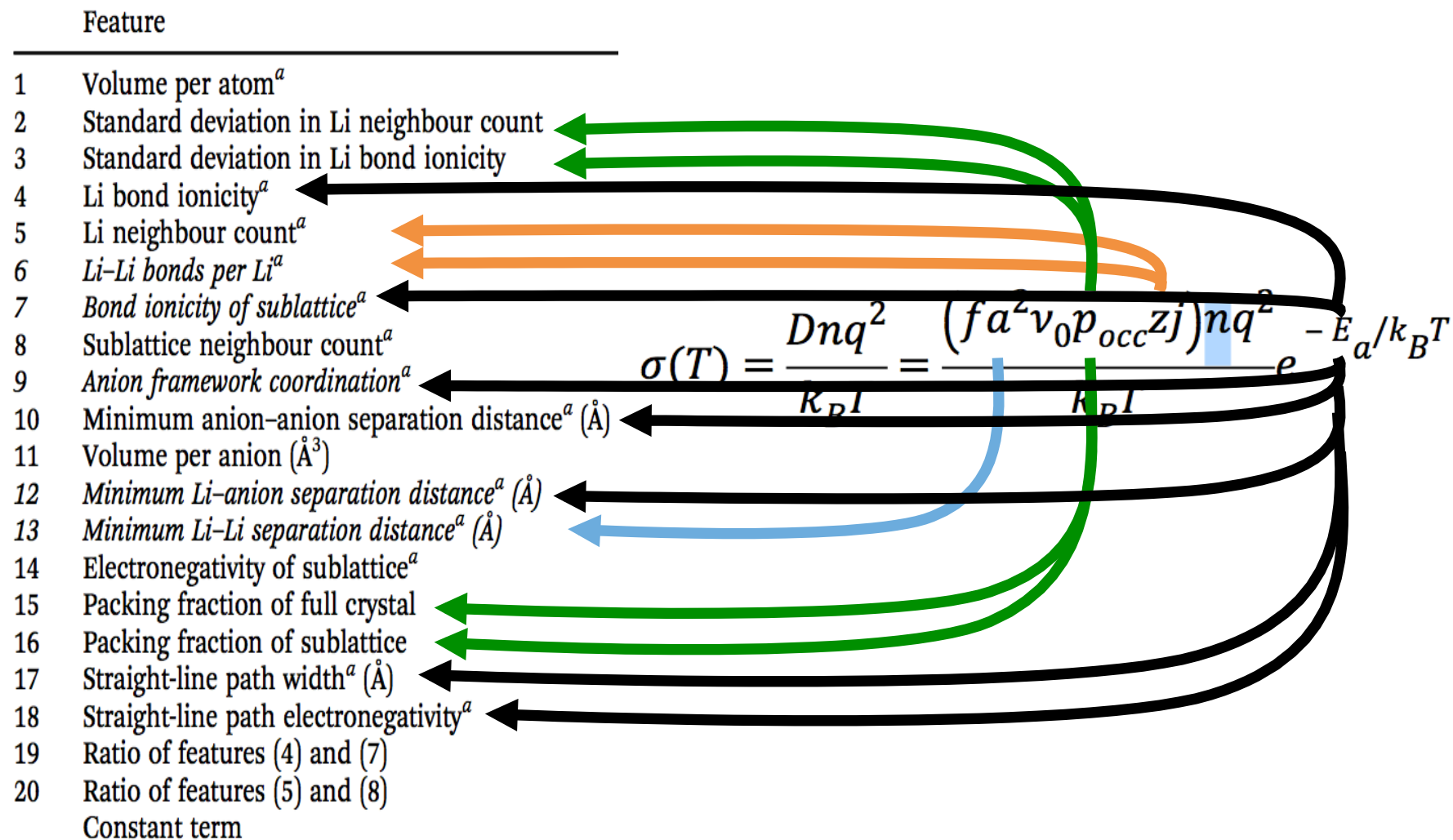
# WE DRAW ON WISDOM/PROPOSALS IN THE LITERATURE FOR READILY COMPUTABLE FEATURES (NO DFT!)

	Feature
1	Volume per atom <sup>a</sup>
2	Standard deviation in Li neighbour count
3	Standard deviation in Li bond ionicity
4	Li bond ionicity <sup>a</sup>
5	Li neighbour count <sup>a</sup>
6	<i>Li-Li bonds per Li<sup>a</sup></i>
7	<i>Bond ionicity of sublattice<sup>a</sup></i>
8	Sublattice neighbour count <sup>a</sup>
9	<i>Anion framework coordination<sup>a</sup></i>
10	Minimum anion-anion separation distance <sup>a</sup> (Å)
11	Volume per anion (Å <sup>3</sup> )
12	<i>Minimum Li-anion separation distance<sup>a</sup> (Å)</i>
13	<i>Minimum Li-Li separation distance<sup>a</sup> (Å)</i>
14	Electronegativity of sublattice <sup>a</sup>
15	Packing fraction of full crystal
16	Packing fraction of sublattice
17	Straight-line path width <sup>a</sup> (Å)
18	Straight-line path electronegativity <sup>a</sup>
19	Ratio of features (4) and (7)
20	Ratio of features (5) and (8)
	Constant term

A physics-based model for a single crystal, with implicit assumptions:

$$\sigma(T) = \frac{Dnq^2}{k_B T} = \frac{(fa^2v_0p_{occ}zj)nq^2}{k_B T} e^{-E_a/k_B T}$$

# WE DRAW ON WISDOM/PROPOSALS IN THE LITERATURE FOR READILY COMPUTABLE FEATURES (NO DFT!)





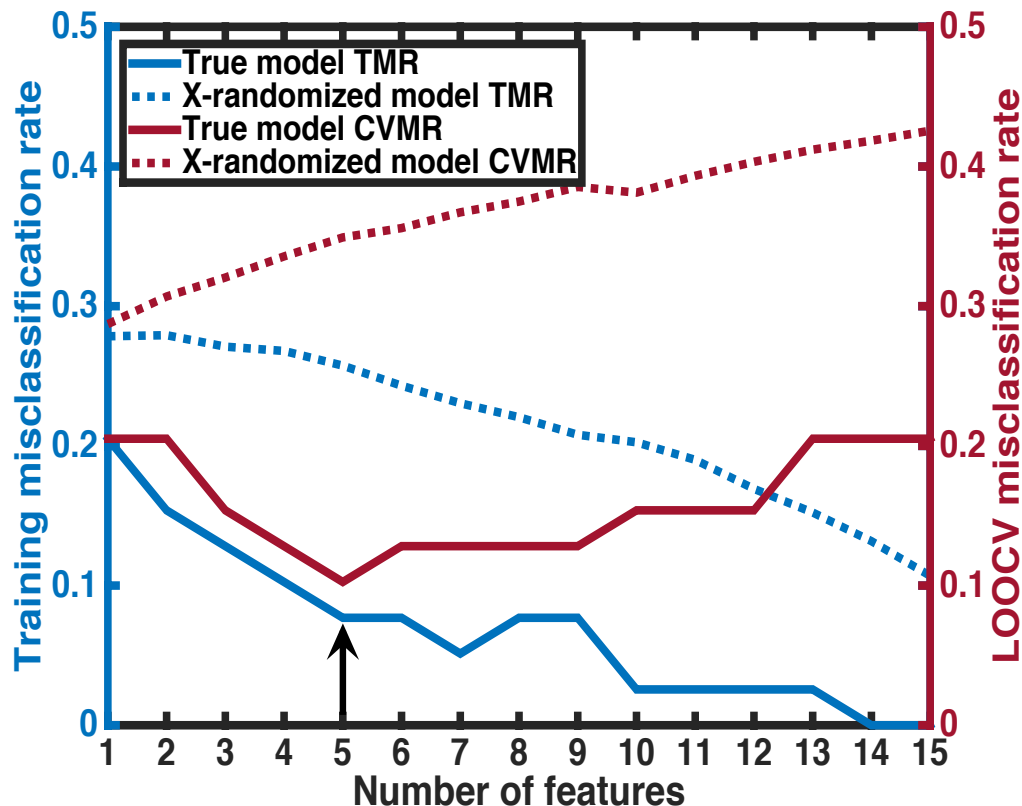
## WE EMPLOY LOGISTIC REGRESSION (TWO-CLASS CLASSIFIER)

Assuming a logistic form, we search for the maximally predictive set of features

$$P_{\text{superionic}}(\mathbf{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \mathbf{x})}$$

$$\boldsymbol{\theta}^T \mathbf{x} = ?$$

# WE EMPLOY LEAVE-ONE OUT CROSS VALIDATION TO DETERMINE OPTIMAL FEATURES

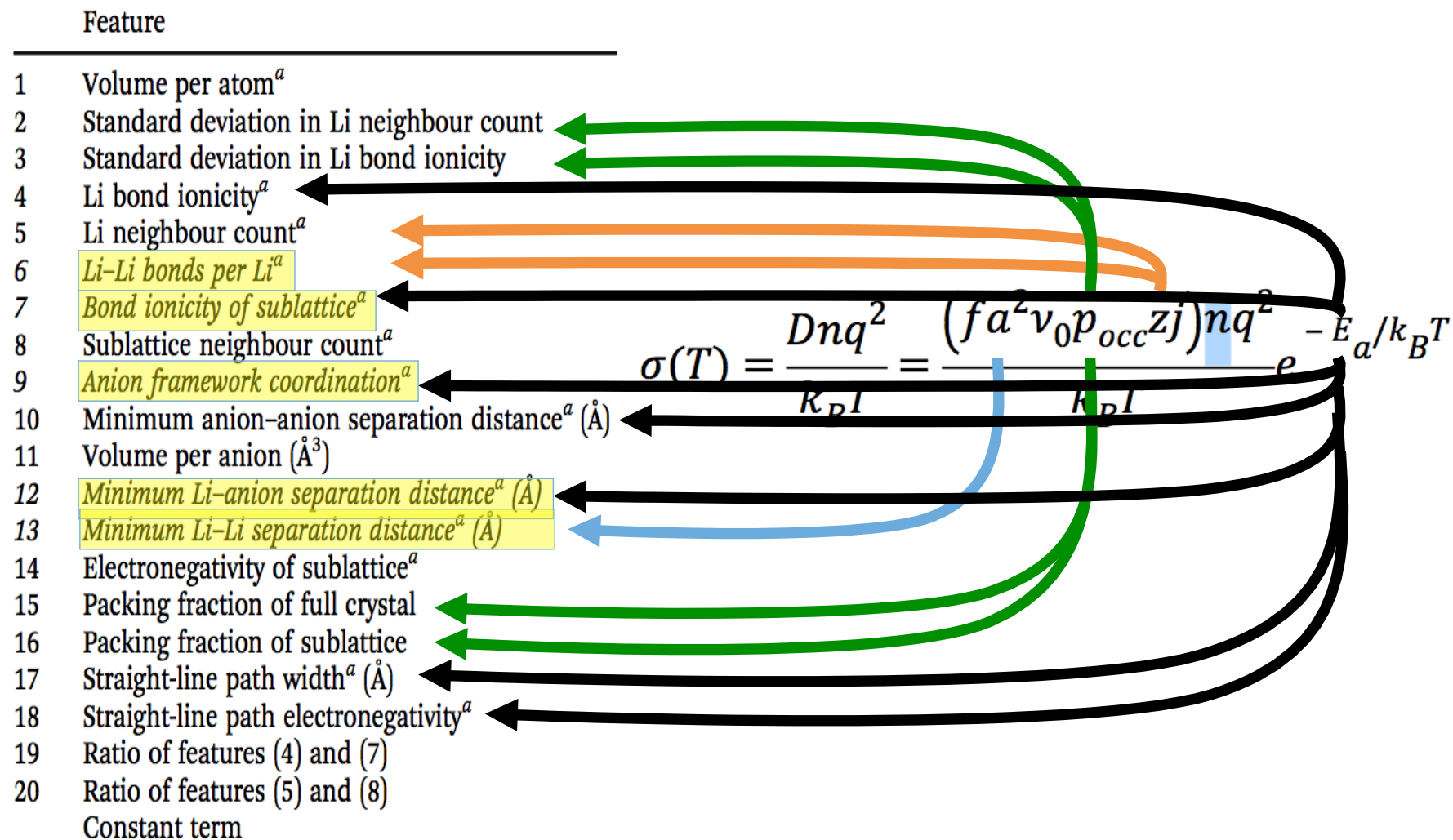


$$\text{TMR} = \frac{1}{M} \sum_{i=1}^M 1 \left\{ \tilde{\sigma}^{(i)} \neq \hat{\sigma}^{(i)} \right\}$$

$$\text{CVMR} = \frac{1}{M} \sum_{i=1}^M 1 \left\{ \tilde{\sigma}^{(i)} \neq \hat{\sigma}_{\text{LOO}}^{(i)} \right\}$$

- We search over all possible combinations of 20 features ( $>10^6$  models)
- Optimal leave-one-out cross-validated misclassification rate = 10%
- Optimal model performance against random guessing: 3-4x improvement

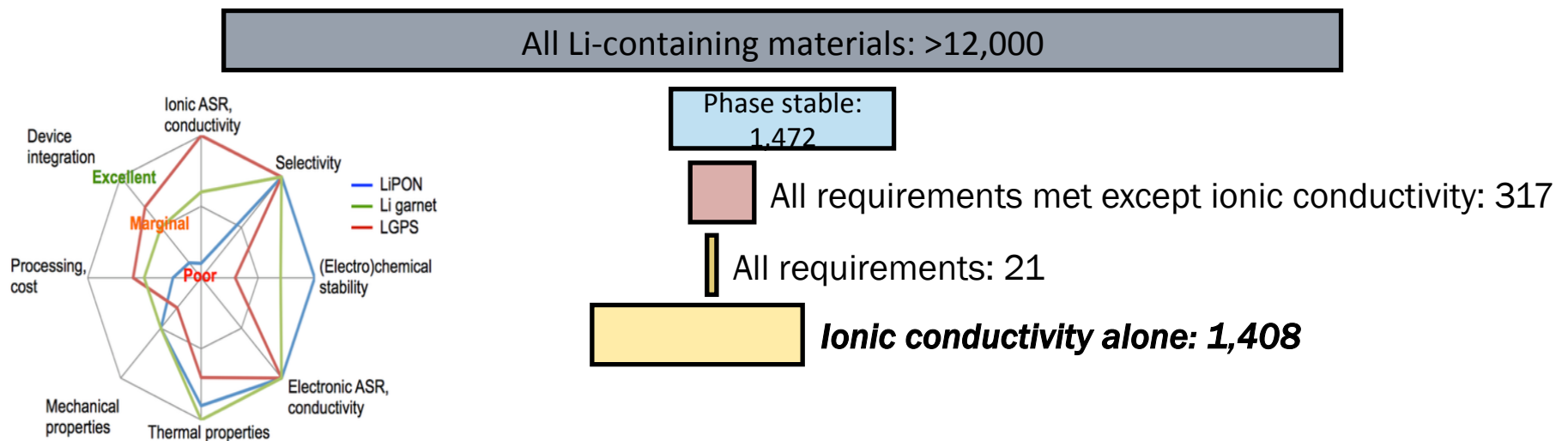
# WE DISCOVER 5 FEATURES THAT BEST CLASSIFY ION CONDUCTORS



# WE PERFORM THE FIRST HOLISTIC STRUCTURE SCREENING OF ALL >12,000 LI-CONTAINING SOLIDS IN THE MATERIALS PROJECT DATABASE

Ionic conductivity is not all that matters! We also screen for:

- High stability against oxidation ← Gibbs free energy of atomization
- High stability against reduction ← Presence of transition metals
- Low electronic conductivity ← Band gap
- High phase stability ← Convex hull
- Low cost ← Cost of raw elements involved
- High earth abundance ← Abundance of elements in Earth's crust



A. D. Sendek, E. J. Reed, et al, Energy and Environmental Science (2017).



# WE PROPOSE 21 NEW PROMISING SOLID ELECTROLYTE CANDIDATES BY SCREENING >12,000

MPID	Chemical formula	$P_{LR}$	$d$	$\epsilon$	$A$	$E_{gap}$	$\tilde{V}_{ox}$	USD/m <sup>2</sup> (10 $\mu$ m thick)	$I_A$	Related study
mp-554076	BaLiBS <sub>3</sub>	0.589	1.049	0.048	1	2.153	9.697	23	38	
mp-532413	Li <sub>5</sub> B <sub>7</sub> S <sub>13</sub>	0.897	1.228	0.024	1	3.553	5.454	42	38	95
mp-569782 <sup>a</sup>	Sr <sub>2</sub> LiCBr <sub>3</sub> N <sub>2</sub>	1.000	6.852	0.000	0	3.973	13.968	16	45	
mp-558219	SrLi(BS <sub>2</sub> ) <sub>3</sub>	0.518	1.556	0.114	1	2.91	13.964	38	38	
mp-15797	LiErSe <sub>2</sub>	0.543	1.505	0.056	1	1.615	6.778	170	67	
mp-29410	Li <sub>2</sub> B <sub>2</sub> S <sub>5</sub>	0.994	1.855	0.003	1	2.538	4.895	29	38	95
mp-676361	Li <sub>3</sub> ErCl <sub>6</sub>	0.655	0.974	0.042	1	5.211	7.794	70	44	96 and 97
mp-643069 <sup>a</sup>	Li <sub>2</sub> HfO	0.652	2.081	0.079	0	4.319	4.054	2.40	60	
mp-19896	Li <sub>2</sub> GePbS <sub>4</sub>	0.604	1.063	0.090	1	2.265	4.591	13	54	90
mp-7744 <sup>a</sup>	LiSO <sub>3</sub> F	1.000	4.097	0.000	0	5.792	13.446	10	34	
mp-22905 <sup>b</sup>	LiCl	0.837	1.381	0.031	1	6.25	4.214	0.94	34	98
mp-34477	LiSmS <sub>2</sub>	0.89	1.33	0.028	1	1.921	8.536	6.50	40	
mp-676109	Li <sub>3</sub> InCl <sub>6</sub>	0.656	1.013	0.058	1	3.373	6.215	5.50	63	96 and 97
mp-559238	CsLi <sub>2</sub> BS <sub>3</sub>	0.812	1.642	0.055	1	3.094	4.798	160	49	
mp-866665 <sup>a</sup>	LiMgB <sub>3</sub> (H <sub>9</sub> N) <sub>2</sub>	1.000	5.149	0.000	0	6.511	11.222	30	38	
mp-8751	RbLiS	0.775	1.279	0.051	1	2.745	4.22	240	34	
mp-15789	LiDyS <sub>2</sub>	0.901	1.339	0.025	1	1.935	8.736	9.20	39	
mp-15790	LiHoS <sub>2</sub>	0.899	1.327	0.025	1	1.965	8.749	300	55	
mp-15791	LiErS <sub>2</sub>	0.899	1.319	0.025	1	2.008	8.761	190	44	
mp-561095 <sup>a</sup>	LiHo <sub>3</sub> Ge <sub>2</sub> (O <sub>4</sub> F) <sub>2</sub>	0.984	3.247	0.009	0	4.163	53.18	370	55	
mp-8430	KLiS	0.76	1.243	0.052	1	3.057	4.348	14	34	

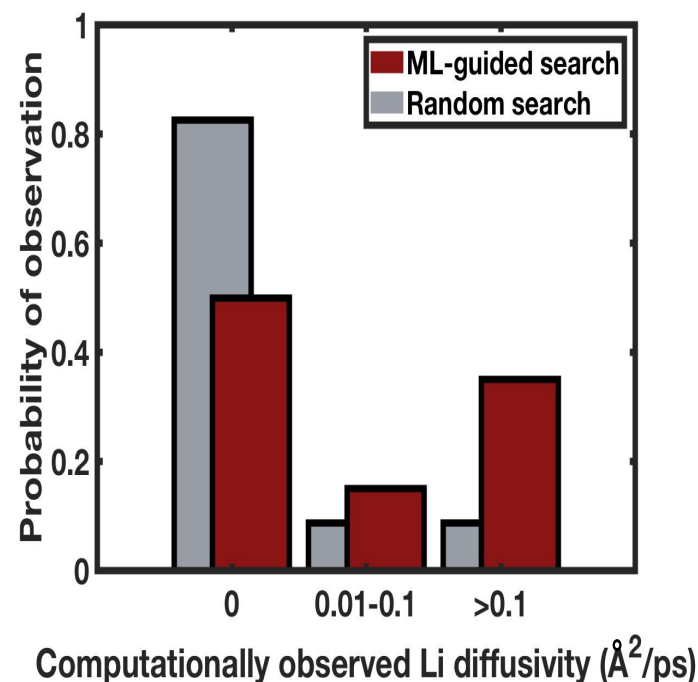
A. D. Sendek, E. J. Reed, et al, Energy and Environmental Science (2017).

How well does the model do?

# We discover ten new solids that are superionic conductors, doubling known superionic conductors

Model success  
rate: 39%  
(melting cases  
excluded)

Random success  
rate: 14%  
(melting cases  
excluded)



Learning from only 40 data points = 3x improvement over guesswork

# TWO METRICS FOR MODEL PERFORMANCE

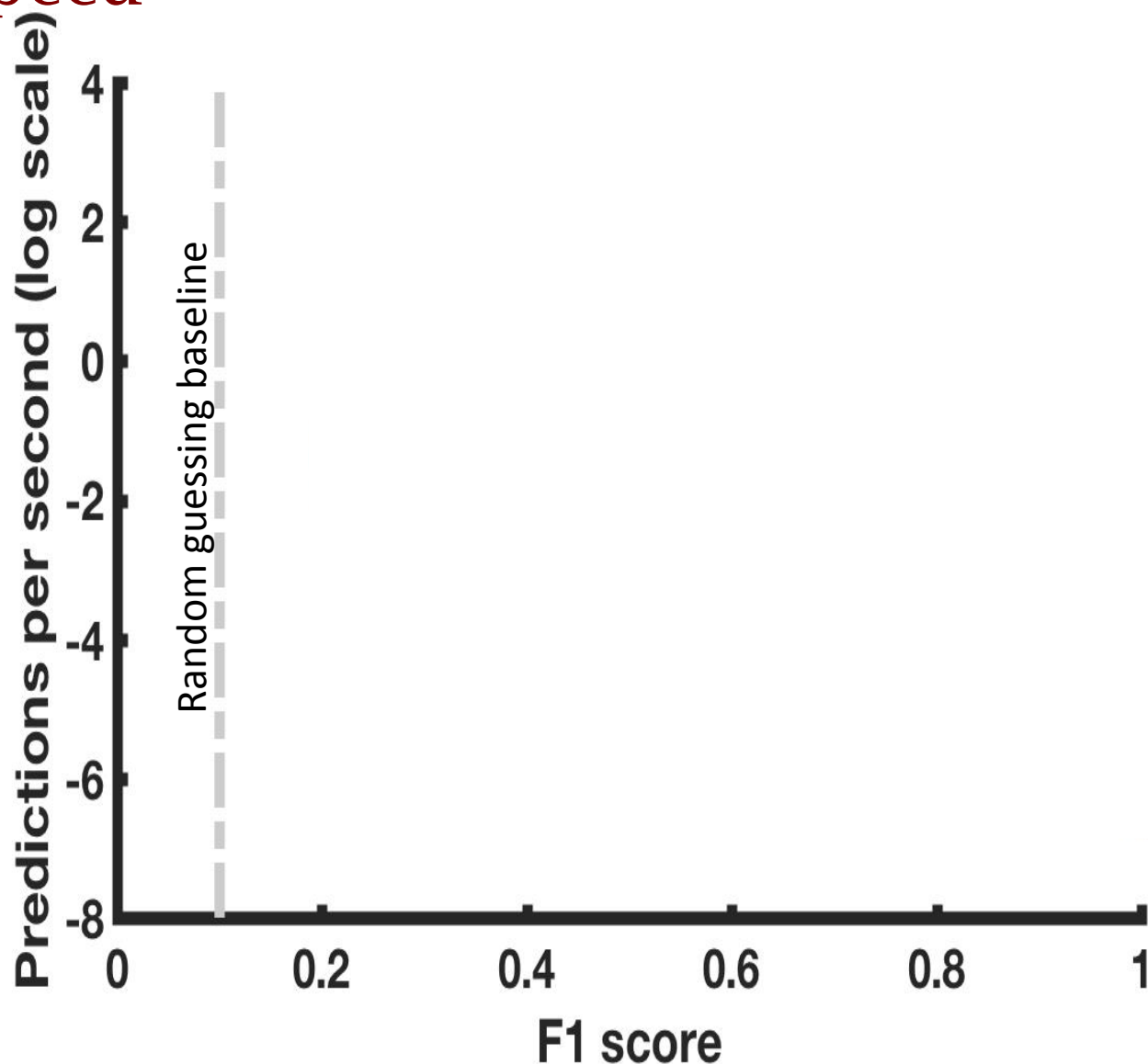


- How much better than random guessing is our model?
  - Approximately three times better
  - Quantifies the state of humanity's knowledge
- Are the false positives acceptable in number?
  - 61% are false positives
  - Need to synthesize two materials to get one that works
  - This is probably good enough in practice, but could be better

Scientists aren't truly guessing at random -  
but can they beat the machine?

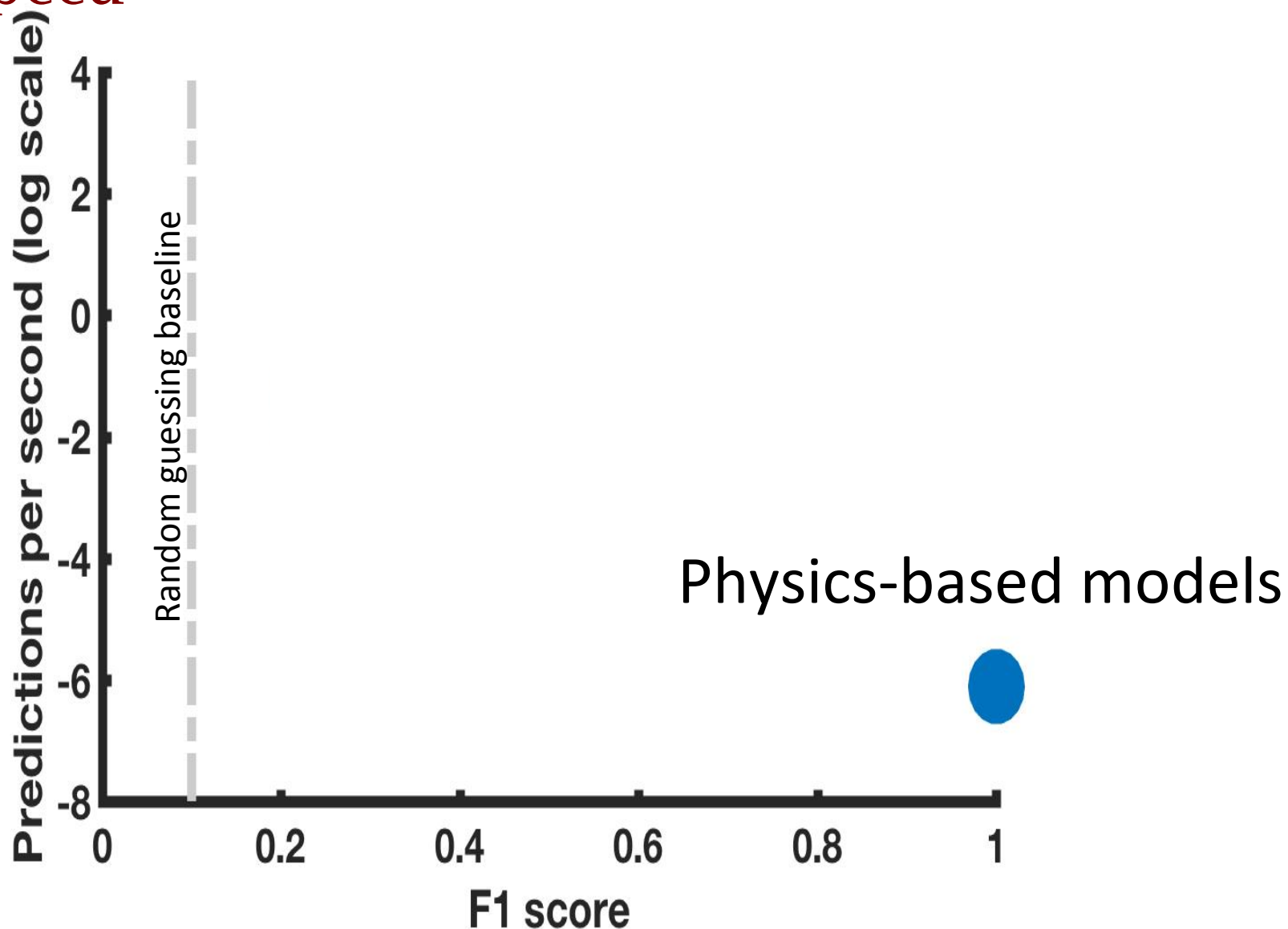


Our algorithm outperforms humans in accuracy and speed



AD Sendek, ED Cubuk, G Cheon, ER Antoniuk, Y Cui, E. J. Reed, Chemistry of Materials (2018).

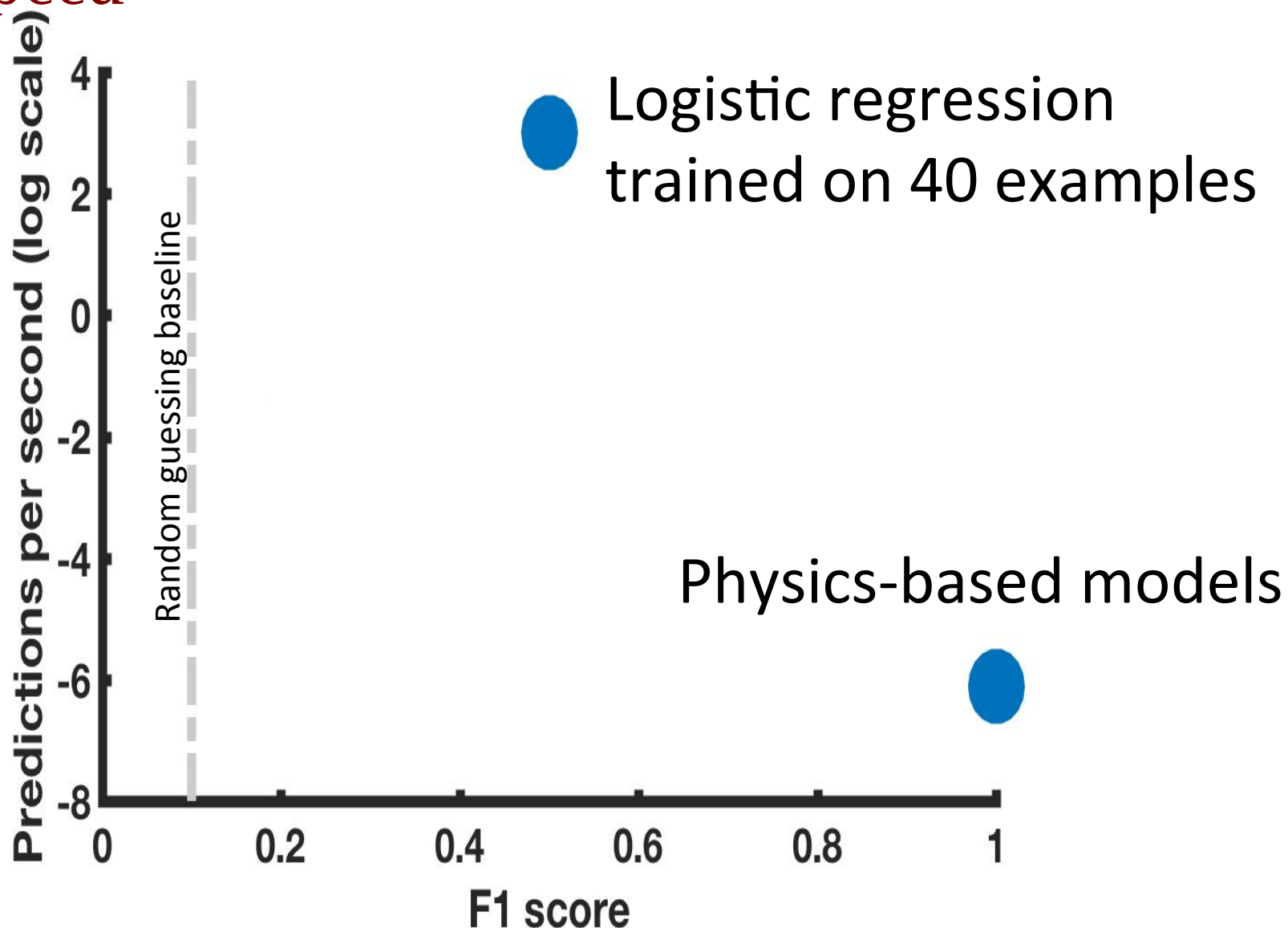
Our algorithm outperforms humans in accuracy and speed



AD Sendek, ED Cubuk, G Cheon, ER Antoniuk, Y Cui, E. J. Reed, Chemistry of Materials (2018).

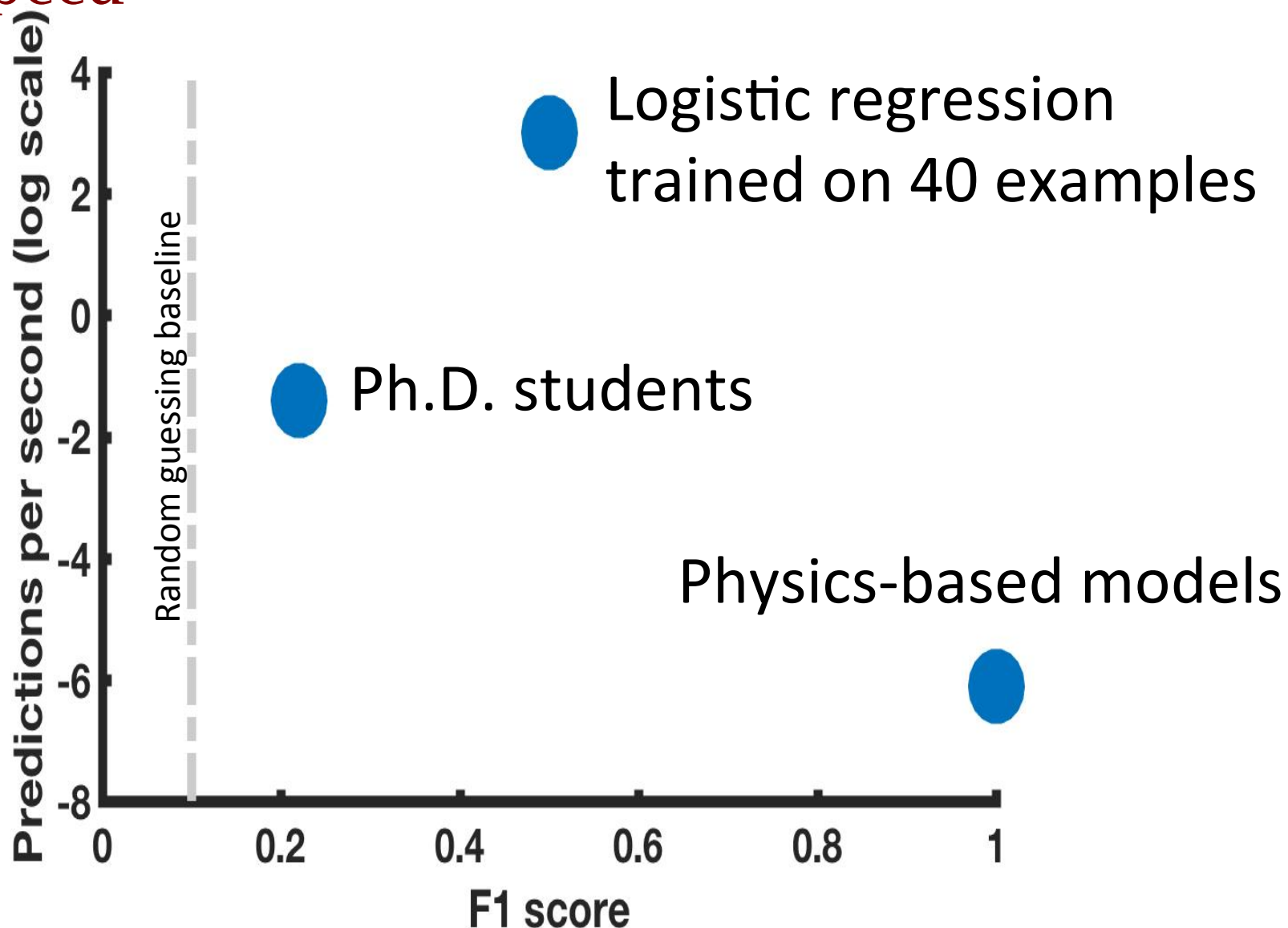


Our algorithm outperforms humans in accuracy and speed



AD Sendek, ED Cubuk, G Cheon, ER Antoniuk, Y Cui, E. J. Reed, Chemistry of Materials (2018).

Our algorithm outperforms humans in accuracy and speed



AD Sendek, ED Cubuk, G Cheon, ER Antoniuk, Y Cui, E. J. Reed, Chemistry of Materials (2018).

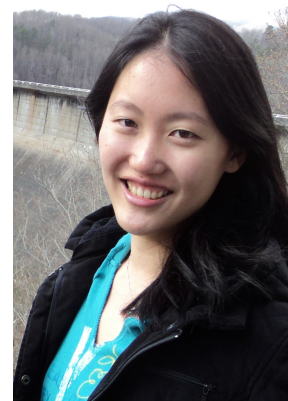
# ACKNOWLEDGEMENTS



Austin Sendek  
(now CEO AIONICS)



Ekin Dogus Cubuk  
(now at Google  
Brain)



Qian Yang  
(now at U  
Conn)



Gowoon Cheon



This work is supported by:

- Stanford TomKat Center for Sustainable Energy Seed grant
- Stanford Graduate Fellowship program

# Acknowledgements

This work is supported by:

- Army High Performance Computing Research Center(AHPCRC)
- Army Research Office W911NF-15-1-0570
- NSF EECS-1436626 and DMR-1455050
- Office of Naval Research N00014-15-1-2697

G.Cheon et al., Nano Letters, 2017  
Data Mining for New Two- and One-Dimensional Weakly  
Bonded Solids and Lattice-Commensurate Heterostructures



Gwoon Cheon



Also:

- Austin Sendek
- Karel-Alexander Duerloo
- Chase Porter
- Yuan Chen