

## ( Big ) Data of Materials Science from First Principles -- Critical Role of the Descriptor --

Matthias Scheffler (\*)

Fritz-Haber-Institut der Max-Planck-Gesellschaft, Berlin; <http://th.fhi-berlin.mpg.de/>

From *the periodic table of the elements* to *a chart of materials*:

Organize materials according to their properties and functions

- figure of merit of thermoelectrics (as function of  $T$ )
- turn-over frequency of catalytic materials (as function of  $T$  and  $p$ )
- efficiency of photovoltaic systems
- etc.



Dmitri Mendeleev  
(1834-1907)

PERIODIC TABLE OF THE ELEMENTS



(\*) Work performed in collaboration with **Luca Ghiringhelli**, **Jan Vybiral**, **Claudia Draxl**, Sergey Levchenko, Alexandre Tkatchenko, Patrick Rinke, Xinguo Ren, and Igor Ying Zhang

## Materials Genome Initiative for Global Competitiveness



Compute the basic properties („genes“) of many (ten or hundred thousand) materials and disseminate that information to the materials community to enable rapid searches of materials properties and help design improved materials.

To help business discover, develop, and deploy new materials twice as fast, we're launching what we call the Materials Genome Initiative. The invention of silicon circuits and lithium ion batteries made computers and iPods and iPads possible, but it took years to get those technologies from the drawing boards to the market place. We can do it faster.

President Obama

Carnegie Mellon University, June 2011

“twice as fast, at a fraction of the cost”



Materials Genome Initiative for Global Competitiveness [http://www.whitehouse.gov/sites/default/files/microsites/stp/materials\\_genome\\_initiative-final.pdf](http://www.whitehouse.gov/sites/default/files/microsites/stp/materials_genome_initiative-final.pdf)

## The Four V of Big Data and an A

Data – data – data (analog to Moore's law)  
numbers, arrays, figures, movies, ...

(so far: most data are not used and even thrown away)



Big-Data Challenge: "four V":

**Volume** (amount of data),

**Variety** (heterogeneity of form and meaning of data),

**Veracity** (uncertainty of data quality),

**Velocity** at which data may change or new data arrive.

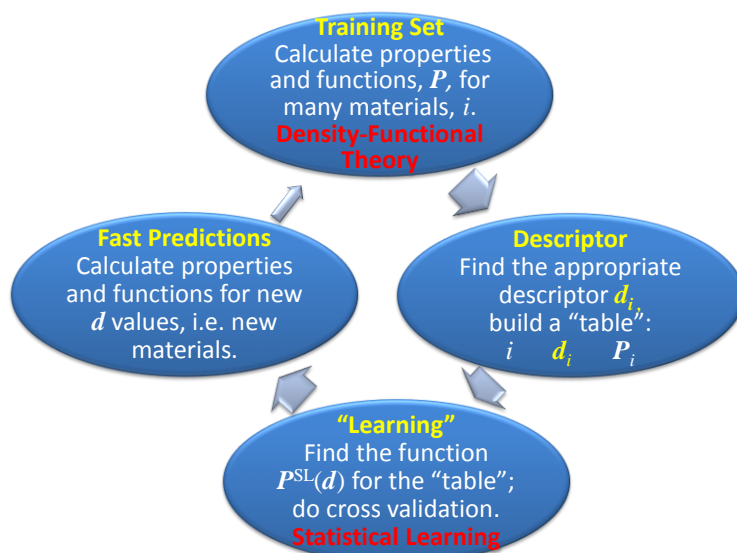
Computed data: Query and read out what was stored. (high-throughput screening)  
**Shouldn't we do more?!**



The four V should be complemented by an "A", **Big-Data Analytics**:

- identify (so far) hidden trends,
- which materials should be studied next as most promising candidates,
- identify anomalies,
- identify the mechanisms that govern a certain material property or function.

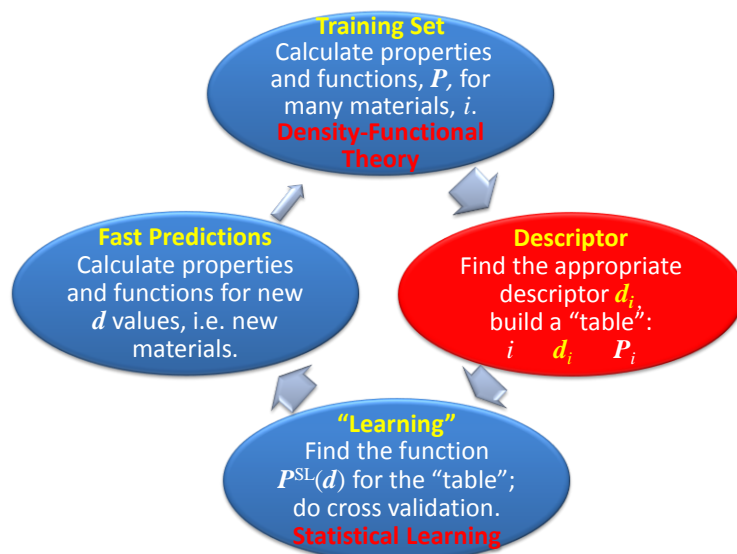
## Big-Data Analytics: How to Arrange the Data



$\{Z, N_f\}$ ,  $T$ ,  $\{p\}$  determine the many-body hamiltonian and statistical mechanics

Statistical mechanics does not tell us what the relevant variables are. This is our choice. If we choose well, the results may be useful, if we chose badly, the results (while formally correct) will probably be useless. (Robert Zwanzig)

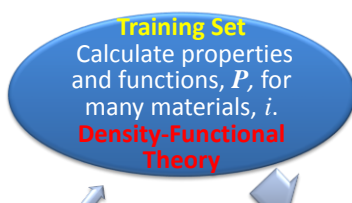
## Big-Data Analytics: How to Arrange the Data



$\{Z, N_f\}, T, \{p\}$  determine the many-body hamiltonian and statistical mechanics

Statistical mechanics does not tell us what the relevant variables are. This is our choice. If we choose well, the results may be useful, if we chose badly, the results (while formally correct) will probably be useless. (Robert Zwanzig)

## Big-Data Analytics: How to Arrange the Data



$\{Z, N_f\}, T, \{p\}$  determine the many-body hamiltonian and statistical mechanics

$d$  characterizes the relevant mechanisms that govern the observed property/function  $P$ . The  $d \rightarrow P^{\text{SL}}$  mapping is complex; identifying the descriptor  $d$  from known data  $P_i$ , is an ill-posed problem (statistical-learning theory): **A little error in the data  $P_i$  may suggest a different descriptor  $d$ . Thus, knowledge of the accuracy of data  $P_i$  is crucial (veracity).** The choice of  $d$  is not unique.

**A) Veracity:** Accuracy of state-of-the-art density-functional theory (validation and verification)

**B) Descriptor:** How to find it, how to understand the causality between  $d$  and  $P^{\text{SL}}$ ?

http://nomad-repository.eu/cms/ Welcome to the NoMaD Rep...


**The NoMaD Repository**

by C. Draxl, E. Blokin, L. Ghiringhelli, F. Mohamed, Ch. Carbogno, M. Scheffler

*Insight by sharing*

Home NoMaD Team Why sharing? Upload Download DOIs Terms and conditions NoMaD repository Other repositories

## Welcome to the NoMaD Repository



**The NoMaD (Novel Materials Discovery) Repository** was established to host, organize, and share materials data.

**NoMaD** copes with the increasing demand and requirement of storing scientific data and making them available for longer periods. Rules of good scientific practice set by many funding agencies, worldwide, require keeping scientific data for 10 years. **NoMaD** offers this for free. **NoMaD** also facilitates research groups to share and exchange


### News

Currently, the **NoMaD Repository** contains **89,464** entries.

Feb. 6, 2015  
White House (USA):  
"It's Time to Open Materials Science Data" ([link](#))

[Open positions](#)

[DOIs](#) for datasets can be




scientific practice set by many funding agencies, worldwide, require keeping scientific data for 10 years. **NoMaD** offers this for free. **NoMaD** also facilitates research groups to share and exchange their results, inside a single group or between two or more, and to recall what was actually done some years ago.

The **NoMaD Repository** enables the confirmatory analysis of materials data, their reuse, and repurposing.

Upload of data is possible without any barrier. Results are accepted in their raw format as produced by the underlying code. The only condition is that the list of authors is provided, and code and code version can be retrieved from the uploaded files. These data can be restricted to the owner or made available to other people (selected by the owner). They can be updated and downloaded at any time.

Read more details concerning the [upload](#). Please, [register](#) or [login](#) to participate.

At present, the repository contains *ab initio* electronic-structure data from density-functional theory and methods beyond. At a later stage, it will be extended by force-field studies and by experimental data. We also give an [outlook on the NoMaD Laboratory](#) that will be dedicated to a *Materials Encyclopaedia*, as the basis for complex queries and the development of various data-analytics tools.



**Codes:** Abinit, crystal, exciting, CASTEP, FHI-aims, Quantum Espresso, VASP – more coming; various xc functionals

[Science Data](#) ([link](#))

[Open positions](#)

[DOIs](#) for datasets can be requested

[Check](#) for related **conferences and workshops**.

We are moving to the HPC Center of the Max Planck Society (RZG). We apologize for any possible instability during the next 2 days.

The **NoMaD Repository** is about joining [eudat](#).








[Financial Support](#)

## Veracity – Validation and Verification

- accuracy of materials-science codes: basis sets, relativity, pseudopotentials, other numerical approximations (verification)
- accuracy of the exchange-correlation functional (validation)

## Veracity – Validation and Verification

### Comparing Solid State DFT Codes, Basis Sets and Potentials

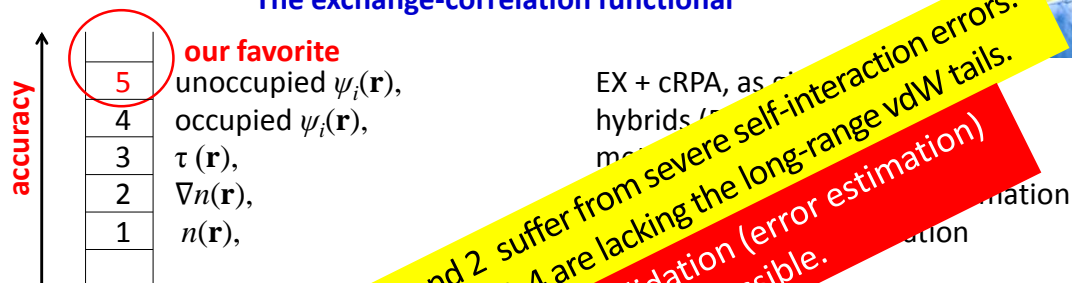
Code	Version	Basis	Electron treatment	$\Delta$ -value	Authors
WIEN2k 	13.1	LAPW/APW+lo	all-electron	0 meV/atom	S. Cottenier
FHI-aims 	081213	tier2 numerical orbitals	all-electron (relativistic atomic_zora scalar)	0.2 meV/atom	ASE [2]
Exciting 	development version	LAPW+xlo	all-electron	0.2 meV/atom	Exciting [10]
FHI-aims 	081213	tier2 numerical orbitals	all-electron (relativistic zora scalar 1e-12)	0.4 meV/atom	ASE [2]
CASTEP 	8.0	plane waves	OTFG CASTEP 8.0	0.5 meV/atom	CASTEP [7]
ABINIT 	7.7.3	plane waves	PAW JTH v0.2 	0.6 meV/atom	F. Jollet and M. Torrent

K. Lejaeghere, V. Van Speybroeck, G. Van Oost, and S. Cottenier, Crit. Rev. Solid State Mater. Sci. 39, 1-24 (2014); <https://molmod.ugent.be/deltacodesdft>. Reference code: WIEN2k

# Veracity – Approximate Treatment of Exchange-Correlation

## Perdew's Dream: A Jacob's Ladder

### The exchange-correlation functional



$\tau(\mathbf{r})$ : Kohn-Sham

EX: exact

cRPA: correlation

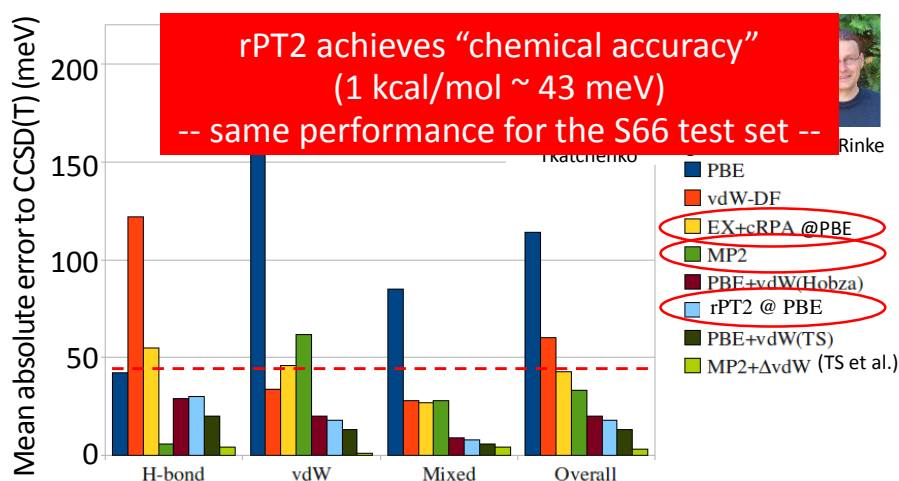
ACFD: adiabatic connection fluctuation-dissipation theorem

Bohm, Pines (1953); Gell-Mann, Brueckner (1957);

Gunnarsson, Lundqvist (1975, 1976); Langreth, Perdew (1977);

X. Ren, P. Rinke, C. Joas, and M. S., *Invited Review, Mater. Sci.* 47, 21 (2012)

## Performance of rPT2 for Weak Intermolecular Interactions: S22 Test Set



CCSD(T): Jurecka, Sponer, Cerny, Hobza, *PCCP* (2006). Langreth-Lundqvist: Gulans, Puska, Nieminen, *PRB* (2009);

rPT2: X. Ren et al. *PRL* (2011) and *NJP* (2013). TS: A. Tkatchenko and M.S., *PRL* (2009); A. Tkatchenko et al., *JCP* (2009)



## Test Sets for Materials Science and Engineering?

Chemists have shown the way. For small and light molecules they developed test sets: G2, NHTBH38, HTBH38, S22, S66 ...

We need a materials test set! We can now do renormalized second-order perturbation theory (similar to CCSD) and even full CI<sup>(\*)</sup> – for certain systems.

Comparison with experiment is very important as well (adsorption energies of molecules, *e.g.* by microcalometry). However, theory-theory comparison is better defined.

(\*) G. H. Booth, A. J. W. Thom, and A. Alavi, J. Chem. Phys. 131, 054106 (2009).  
G. H. Booth, A. Grüneis, G. Kresse, and A. Alavi, Nature 493, 365 (2013).

## Test set for materials science and engineering




7 elements and 12 binaries  
with cubic structure  
(for the start)

H	Light main group elements						He
Li	Be	B	C	N	O	F	Ne
Na	Mg	Al	Si	P	S	Cl	Ar
K	Ca	Ga	Ge	As	Se	Br	Kr
Rb	Sr	In	Sn	Sb	Te	I	Xe
Cs	Ba						

Ne, Ar, Al (fcc); Li, Na (bcc); C, Si (diamond);  
LiH, LiF, LiCl, NaF, NaCl, MgO, MgS (rocksalt);  
BeS, BP, AlP, SiC, BN (zincblende)



- **MSE properties:** cohesive, electronic, elastic and vibrational
- **Representative** for cubic metals, semiconductors, and insulators
- **Numerically accurate reference values from theory**,  
incl. MP2, RPA, CCSD(T)



ABOUT  
SEARCH  
LINKS  
CITE

http://mse.fhi-berlin.mpg.de

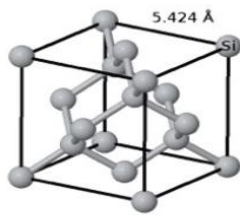
**TEST SET FOR MATERIALS SCIENCE AND ENGINEERING**

Group	Material	Structure	Method	$E_{\text{coh}}$ (eV)	$a_0$ (Å)	$B$ (GPa)	CS	BS
14	Si	diamond	LDA	5.325	5.402	86.1		
14	Si	diamond	PBE	4.585	5.471	89.1		
14	Si	diamond	PBE+vdW(TS)	4.868	5.448	91.4		
14	Si	diamond	PBE+vdW(MBD)	4.844	5.434	93.4		
14	Si	diamond	PBEsol	4.972	5.434	94.2		
14	Si	diamond	HSE06	4.798	5.424	99.1		
14	Si	diamond	CCSD(T)	4.5				

---

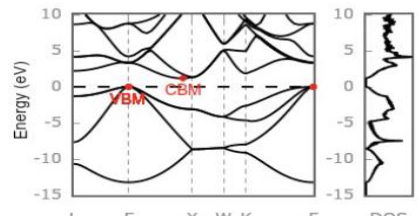
**VISUALIZATION**

Diamond structure, conventional cell



Cell angles:  $\alpha = 90.0^\circ$   
Cell volume = 159.573 Å<sup>3</sup>

Band structure and density of states of Si in diamond structure, lattice parameters optimized with HSE06, band structure calculated with HSE06.



Move the mouse over the bands to see their energies  
Show VBM and CBM

## The Four V of Big Data and an A

Data – data – data (analog to Moore's law)  
numbers, arrays, figures, movies, ...

Computed data: Query and read out what  
was stored. (high-throughput screening)  
*Shouldn't we do more?!*

### Big-Data Analytics

- Finding  $d$  from  $P$
- Causality: the science behind  $P(d)$

Big-Data Challenge: "four V":

*Volume* (amount of data),

*Variety* (heterogeneity of form and meaning of data),

*Veracity* (uncertainty of data quality),

*Velocity* at which data may change or new data arrive.

The four V should be complemented by an "A", the **Big-Data Analytics**:

- identify (so far) hidden correlations,
- which materials should be studied next as most promising candidates,
- identify anomalies,
- identify the mechanisms that govern a certain material property or function.



## Kernel Regression

We have data  $\{P_i\}$  at “coordinates”  $\{x_i\}$       $x_i$  = set of descriptive parameters (descriptor)

$$P_i = P(x_i) = \sum_{k=1}^N c_k K(x_i, x_k)$$

Linear regression:      $K(x_i, x_k) = x_i \cdot x_k$       $P(x_i) = x_i \cdot c^*$

Polynomial kernel      $K(x_i, x_k) = (x_i \cdot x_k + c)^d$

Gaussian kernel      $K(x_i, x_k) = \exp(-\sum_j (x_i - x_k)^2 / 2\sigma_j^2)$

More data means better representation.

Do we “learn” anything?

For successful learning, we need a “good” descriptor:  $P(x_i) \rightarrow P(d_i)$

## Toy Model: Descriptor for the Classification “Zincblende/Wurtzite or Rocksalt?”

Only DFT-

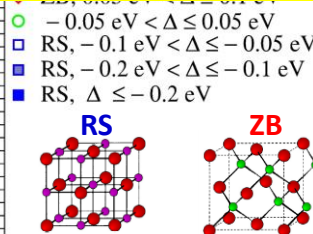
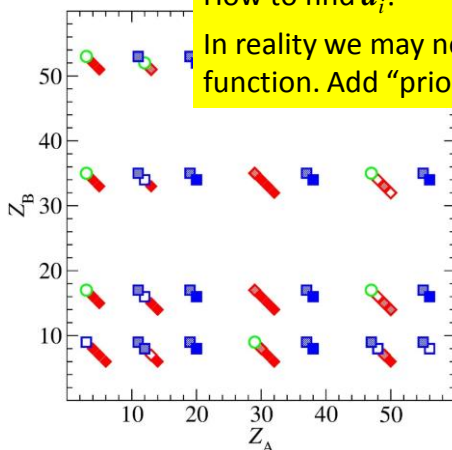
Arrange data  $P(x_i) \rightarrow P(d_i)$ , so that  $P(d_i)$  is a “well behaved function”. Fit this function by machine learning.

$Z_B$ ?

82 oct

How to find  $d_i$ ?

In reality we may not have enough data to learn a complex function. Add “prior knowledge” (prejudice).



between are very small. For Si: 0.01% of the energy of a Si atom, or 0.1% of the 4 valence electrons. Complexity:  $T_s[n]$  and  $E_{xc}$ .

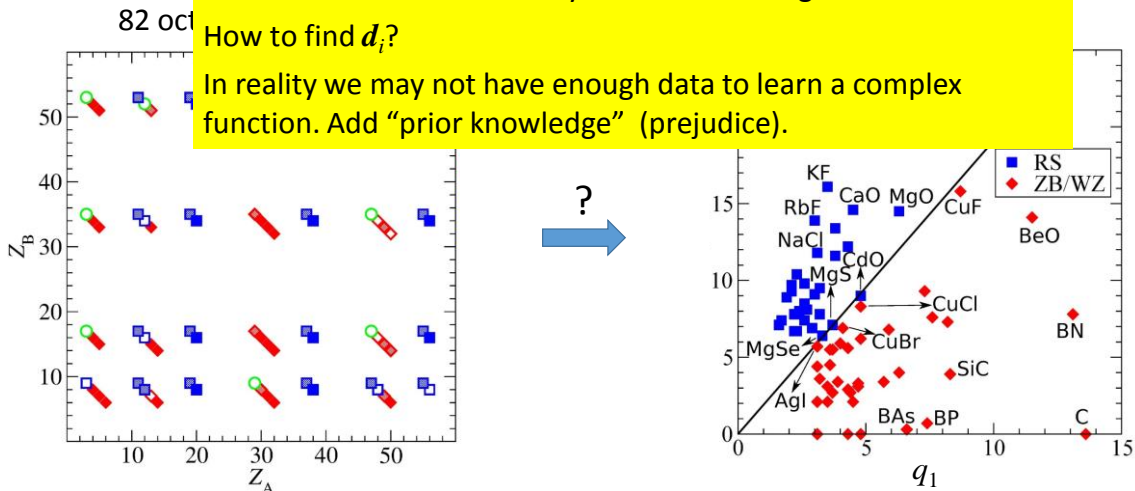
Machine learning can fit the  $P(Z_A, Z_B)$  data well, but fails completely in predictions.

## Toy Model: Descriptor for the Classification “Zincblende/Wurtzite or Rocksalt?”

Arrange data  $P(x_i) \rightarrow P(d_i)$ , so that  $P(d_i)$  is a “well behaved function”. Fit this function by machine learning.

How to find  $d_i$ ?

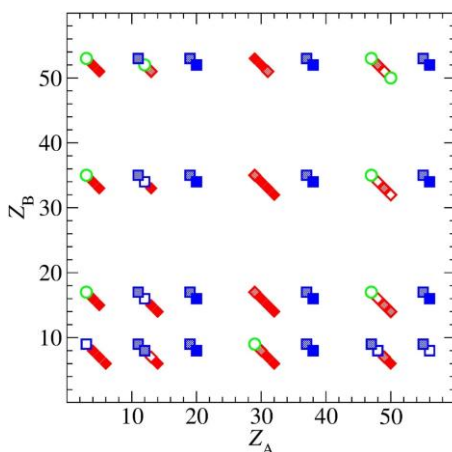
In reality we may not have enough data to learn a complex function. Add “prior knowledge” (prejudice).



## Toy Model: Descriptor for the Classification “Zincblende/Wurtzite or Rocksalt?”

Only DFT-LDA: Can we predict not yet calculated LDA structures from  $Z_A$  and  $Z_B$ ?

82 octet AB binary compounds



### Wish List for a Descriptor

- A descriptor is an array of real numbers that uniquely characterizes the material as well as property-relevant elementary processes.
- Materials that are very different (similar) should be characterized by very different (similar) descriptor values.
- The determination of the descriptor must not involve calculations as intensive as those needed for the evaluation of the property to be predicted.
- The dimension of the descriptor should be as low as possible, but not lower.

## Toy Model: Descriptor for the Classification “Zincblende/Wurtzite or Rocksalt?”

the key scientific challenge: find a descriptor  $d$  that works.

How to find a good descriptor, also for more complicated properties and functions?

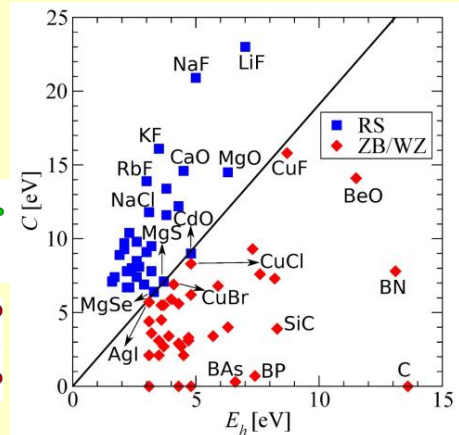
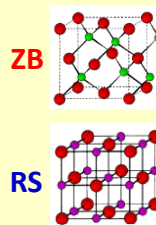
### Example (the traditional way):

J. C. Phillips and J. A. Van Vechten (1969/70)<sup>(\*)</sup> :  
A two-dimensional descriptor that distinguishes  
materials that crystallize in ZB/WZ vs. RS structures:

$E_h =$  } related to crystal's band gap, dielectric  
 $C =$  } constant, nearest-neighbor distance

There are several other descriptors for this  
classification goal, by various authors.

(\*) J. A. Van Vechten, Phys. Rev. B **182**, 891 (1969);  
J. C. Phillips, Rev. Mod. Phys. **42**, 317 (1970)



## Statistical Learning (Machine Learning)

fit and/or interpolation of discrete, known data points  $\{P_i\}$  and building a function  $P(d)$

the key scientific challenge: find a reliable, low dimensional descriptor  $d$ .

### kernel ridge regression

$$P(d) = \sum_{i=1}^N c_i \exp(-\|d_i - d\|_2^2 / 2\sigma^2)$$

$$\sum_{i=1}^N (P(d_i) - P_i)^2 + \text{minimize}$$

$$\lambda \sum_{i,j=1}^{N,N} c_i c_j \exp(-\|d_i - d_j\|_2^2 / 2\sigma^2)$$

$$\|d_i - d_j\|_2^2 = \sum_{\alpha=1}^{\Omega} (d_{i,\alpha} - d_{j,\alpha})^2$$

### linear

R. Tibshirani, J. Royal Statist. Soc. B 58, 267 (1996)

$$P(d) = dc$$

$$\sum_{i=1}^N (P(d_i) - P_i)^2 +$$

$$\lambda \|c\|_1$$

$$\|c\|_1 = \sum_{\alpha=1}^M |c_\alpha|$$

least absolute shrinkage and selection  
operator (LASSO) for feature selection

1) Primary Features, 2) Feature Space, 2) Descriptors

1)

ID	Description	free atoms	Symbols	#
A1	Ionization Potential (IP) and Electron Affinity (EA)		IP(A) EA(A) IP(B) EA(B) [1]	4
A2	Highest occupied (H) and lowest unoccupied (L) Kohn-Sham levels		H(A) L(A) H(B) L(B)	4
A3	Radius at the max. value of <i>s</i> , <i>p</i> , and <i>d</i> valence radial radial probability density		$r_s(A)$ $r_p(A)$ $r_d(A)$ $r_s(B)$ $r_p(B)$ $r_d(B)$	6

ID	Description	free dimers	Symbols	#
A4	Binding energy		$E_b(AA)$ $E_b(BB)$ $E_b(AB)$	3
A5	HOMO-LUMO KS gap		HL(AA) HL(BB) HL(AB)	3
A6	Equilibrium distance		$d(AA)$ $d(BB)$ $d(AB)$	3

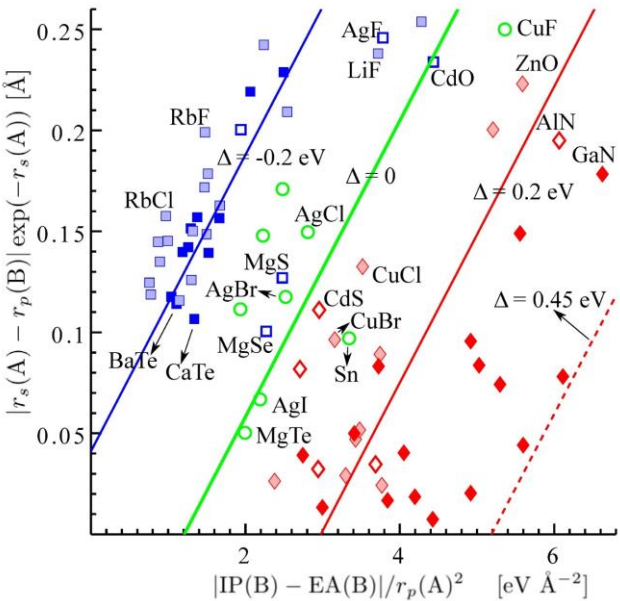
2)

We start with 23 primary features  
and build > 10,000 non linear combinations

3) LASSO finds the most important descriptors:

$$\frac{IP(B) - EA(B)}{r_p(A)^2}, \frac{|r_s(A) - r_p(B)|}{\exp(r_s(A))}, \frac{|r_p(B) - r_s(B)|}{\exp(r_d(A) + r_s(B))}$$

Statistical Learning (Machine Learning): LASSO, 2-Dim. Descriptor



- $\Delta = E(RS) - E(ZB)$
- ◆ ZB,  $\Delta > 0.2$  eV
  - ◆ ZB,  $0.1 \text{ eV} < \Delta \leq 0.2$  eV
  - ◆ ZB,  $0.05 \text{ eV} < \Delta \leq 0.1$  eV
  - $-0.05 \text{ eV} < \Delta \leq 0.05$  eV
  - RS,  $-0.1 \text{ eV} < \Delta \leq -0.05$  eV
  - RS,  $-0.2 \text{ eV} < \Delta \leq -0.1$  eV
  - RS,  $\Delta \leq -0.2$  eV

$P(d) = dc$

The complexity and science is in the  
descriptor (identified from >10,000  
features).

## Statistical Learning (Machine Learning): Descriptor

Mean absolute error (MAE), and maximum absolute error (MaxAE), in eV, (first two lines) and for a leave-10%-out cross validation (CV), averaged over 150 random selections of the training set (last two lines). For  $(Z_A^*, Z_B^*)$ , each atom is identified by a string of three random numbers.

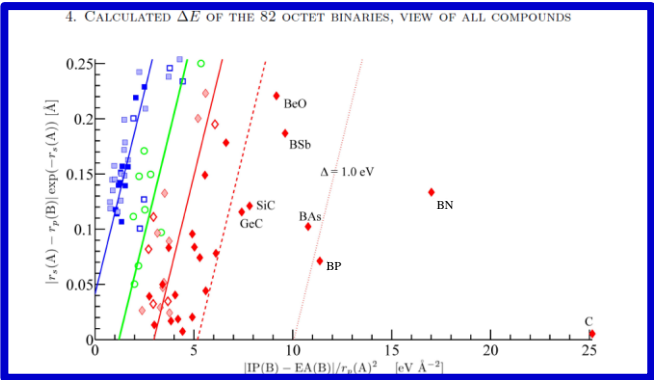
Descriptor	$Z_A, Z_B$	$Z_A^*, Z_B^*$	1D	2D	3D	5D
MAE	$1 \cdot 10^{-4}$	$3 \cdot 10^{-3}$	0.12	0.08	0.07	0.05
MaxAE	$8 \cdot 10^{-4}$	0.03	0.32	0.32	0.24	0.20
MAE, CV	0.13	0.14	0.12	0.09	0.07	0.05
MaxAE, CV	0.43	0.42	0.27	0.18	0.16	0.12

## Statistical Learning (Machine Learning): Descriptor

Mean absolute error (MAE), and maximum absolute error (MaxAE), in eV, (first two lines) and for a leave-10%-out cross validation (CV), averaged over 150 random selections of the training set (last two lines). For  $(Z_A^*, Z_B^*)$ , each atom is identified by a string of three random numbers.

Descriptor	$Z_A, Z_B$	$Z_A^*, Z_B^*$	1D	2D	3D	5D
MAE	$1 \cdot 10^{-4}$	$3 \cdot 10^{-3}$	0.12	0.08	0.07	0.05
MaxAE	$8 \cdot 10^{-4}$	0.03	0.32	0.32	0.24	0.20
MAE, CV	0.13	0.14	0.12	0.09	0.07	0.05
MaxAE, CV	0.43	0.42	0.27	0.18	0.16	0.12

Statistical Learning (Machine Learning): Descriptor



lute error (MaxAE), in eV, (first two CV), averaged over 150 random ( $Z_A^*$ ,  $Z_B^*$ ), each atom is identified by

	1D	2D	3D	5D
MAE, CV	0.12	0.08	0.07	0.05
MaxAE, CV	0.32	0.32	0.24	0.20
	0.12	0.09	0.07	0.05
	0.27	0.18	0.16	0.12

Drawing Causal Inference from Big Data (Scientific Insight)

Correlation between  $d$  and  $P$ , i.e.  $P$  is a function of  $d$ ,  $P(d)$ , reflects causal inference if it is based on sufficient information(\*)

There are four possibilities (types of causality) behind  $P(d)$ :



Judea Pearl

1.  $d \rightarrow P$  :  $P$  “listens” to  $d$
2.  $A \rightarrow d$  and  $A \rightarrow P$  : There is no direct connection between  $d$  and  $P$ , but  $d$  and  $P$  both “listen” to a third “actuator”
3.  $P \rightarrow d$  :  $d$  “listens” to  $P$
4. There is no direct connection between  $d$  and  $P$ , but they have a common effect that listens to both and screams: “I occurred” (Berkson bias; Judea Pearl)

(\*) Construct  $d$  with scientific knowledge (prejudice?), or use “big data” for  $\{P_i\}$ .



## Drawing Causal Inference from Big Data (Scientific Insight)

### Example:

The probability of childhood leukemia is higher for people living close to electricity power lines.

There is no direct connection between leukemia and the electromagnetic field.

Living close to electric power lines is not a desired residence. People living near power lines tend to be poorer than the control group, and there is a relationship between poverty and cancer.

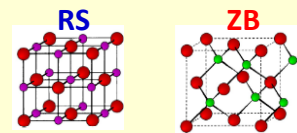
Poverty → higher probability for living close to power lines  
 Poverty → higher chances for cancer

? ? correlation  
 no direct relation  
 causality

## Drawing Causal Inference from Big Data (Scientific Insight)

Our previous example:

**Prediction of the energy difference between ZB/WZ and RS of binary Compound semiconductors**



There is no scientific law that connects the descriptor

$$\frac{IP(B) - EA(B)}{r_p(A)^2}, \frac{|r_s(A) - r_p(B)|}{\exp(r_s(A))}, \frac{|r_p(B) - r_s(B)|}{\exp(r_d(A) + r_s(B))}$$

directly with the total-energy difference (we are not able to write it down).

However,  $Z_A, Z_B$  determine these descriptors,

and  $Z_A, Z_B$  determine the many-body Hamiltonians and the total-energy difference.

Poverty → higher probability for living close to power lines  
 Poverty → higher chances for cancer

? ? correlation  
 no direct relation  
 causality

## Drawing Causal Inference from Big Data (Scientific Insight)

Correlation between  $d$  and  $P$  reflects causal inference  
if it is based on sufficient information<sup>(\*)</sup>

There are four possibilities (types of causality) behind  $P(d)$ :



Judea Pearl

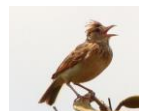
1.  $d \rightarrow P$  :  $P$  "listens" to  $d$
2.  $A \rightarrow d$  and  $A \rightarrow P$  : There is no direct connection between  $d$  and  $P$ , but  $d$  and  $P$  both "listen" to a third "actuator"
3.  $P \rightarrow d$  :  $d$  "listens" to  $P$
4. There is no direct connection between  $d$  and  $P$ , but they have a common effect that listens to both and screams: "I occurred" (Berkson bias; Judea Pearl)

<sup>(\*)</sup> Construct  $d$  with scientific knowledge (prejudice?), or use "big data" for  $\{P_i\}$ .

## Drawing Causal Inference from Big Data (Scientific Insight)

ROMEO: "It was the lark, the bird that sings at dawn, not the nightingale. Look, my love, what are those streaks of light in the clouds parting in the east? Night is over, and day is coming. ... "

case # 3



The *singing of the lark* is a good descriptor for "*the sun will rise soon*".  
The *singing of the lark* is not the actuator of (the mechanism behind) the sunrise.



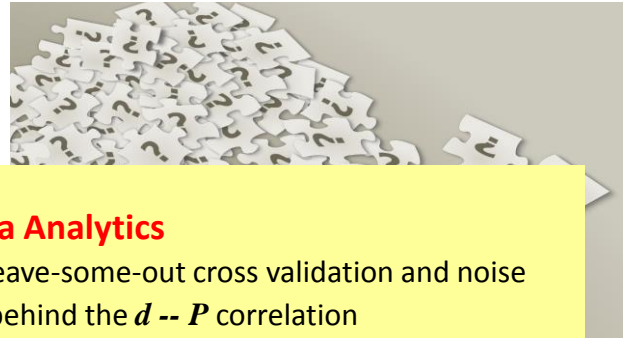
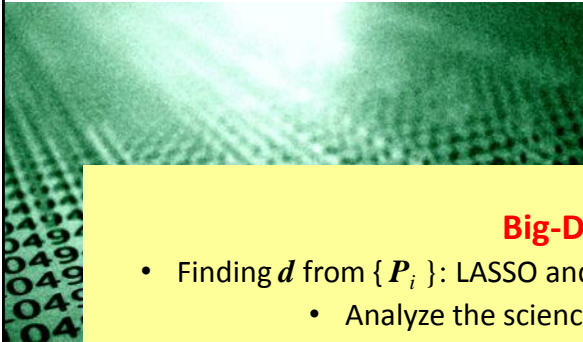
Conclusion / Suggestion: Accept "larks" (not just scientific laws) to predict materials properties.

## The Four V of Big Data and an A



Data – data – data (analog to Moore's law)  
numbers, arrays, figures, movies, ...

Computed data: Query and read out what  
was stored. (high-throughput screening)  
Shouldn't we do more?!



### Big-Data Analytics

- Finding  $d$  from  $\{P_i\}$ : LASSO and leave-some-out cross validation and noise
  - Analyze the science behind the  $d$  --  $P$  correlation
- The big-data challenge in materials science: Look for anomalies, not the crowd

## Summary and Outlook

- Machine learning may find structure in data that is invisible to humans.
- Causal models, i.e. finding causal descriptors, are richer. They are able to provide scientific insight and understanding.
- They can tell how to do machine learning on difficult tasks.
- Question: Why do we want to achieve insight and understanding? Isn't it good enough to have a predictive model?
- Question: Is it possible to assign error bars to predictions (of unexpected situations)?

### Next steps:

- Beyond the linear fit: Non-linear kernel with  $l_1$ -norm regularization.
- Higher accuracy: MaxAE loss function with  $l_1$ -norm regularization.
- Improving the systematic creation of the feature space.